

I10 Patient Segmentation Analysis

Clinical Clustering for Hypertension Management
Using Gower Distance and PAM Algorithm

Metric	Value
Patient Cohort	523 patients
Optimal Clusters	4 segments
Silhouette Score	0.211
Stability (Jaccard)	0.317
Validation Method	Bootstrap (100 iterations)

Executive Summary

This report presents a comprehensive patient segmentation analysis of 523 individuals diagnosed with essential hypertension (ICD-10 code I10). Using advanced clustering techniques that account for mixed clinical, demographic, and utilization data, we identified 4 distinct patient segments with clinically meaningful characteristics.

Study Objectives

- Segment I10 hypertension patients into clinically distinct groups
- Identify patterns in disease severity, comorbidity burden, and healthcare utilization
- Provide actionable insights for targeted care management strategies
- Validate clustering robustness through silhouette analysis and bootstrap stability testing

Key Findings

Segment 1 (202 patients, 38.6%): Median SBP 149 mmHg, Age 56 years, 5 encounters per year.

Segment 2 (65 patients, 12.4%): Median SBP 142 mmHg, Age 56 years, 2 encounters per year.

Segment 3 (124 patients, 23.7%): Median SBP 150 mmHg, Age 66 years, 6 encounters per year.

Segment 4 (132 patients, 25.2%): Median SBP 126 mmHg, Age 67 years, 5 encounters per year.

Clinical Implications

The identified segments demonstrate significant clinical heterogeneity in blood pressure control, comorbidity burden, and healthcare utilization patterns. These findings enable:

- Risk-stratified care pathways tailored to each segment's characteristics
- Resource allocation optimized for high vs. low-utilizer segments
- Targeted interventions addressing specific comorbidity patterns
- Personalized monitoring frequency based on disease severity

1. Introduction

1.1 Background

Essential hypertension (ICD-10 code I10) affects approximately 1.28 billion adults worldwide and remains a leading risk factor for cardiovascular disease, stroke, and chronic kidney disease. Despite the availability of effective antihypertensive therapies, achieving blood pressure control remains challenging due to significant heterogeneity in patient characteristics, comorbidity burden, and treatment response.

Traditional "one-size-fits-all" approaches to hypertension management may not adequately address the diverse needs of the patient population. Patient segmentation offers a data-driven approach to identify clinically meaningful subgroups that share similar characteristics, enabling more precise and personalized care strategies.

1.2 Study Objectives

This analysis aims to segment a cohort of I10 hypertension patients using unsupervised machine learning techniques that can handle mixed clinical, demographic, and utilization data. Specific objectives include:

1. Apply Gower distance-based clustering to accommodate mixed data types (continuous, categorical, binary)
2. Identify the optimal number of clinically distinct patient segments
3. Validate clustering stability and clinical meaningfulness
4. Generate actionable insights for care management and resource allocation

2. Data and Methods

2.1 Study Cohort

The analysis included 523 adult patients with a primary diagnosis of essential hypertension (ICD-10 code I10). Patients were identified from electronic health records and required to have documented blood pressure measurements and demographic information. The cohort represents a diverse population with varying degrees of disease severity, comorbidity burden, and healthcare utilization patterns.

2.2 Feature Engineering

A total of 16 features were selected for clustering, spanning four key domains:

Clinical Severity: Systolic blood pressure (latest), BP stage classification, BMI (latest), BMI class

Demographics: Age, sex

Comorbidity Burden: Total ICD-3 code count, presence of diabetes (E11), dyslipidemia (E78), liver disease (K76), atherosclerosis (I70)

Healthcare Utilization: 12-month encounter count

Data Quality: Missing data indicators for SBP, DBP, and BMI

Missing Data Handling: Numeric features with missing values were imputed using median substitution, while categorical features were imputed using mode values. This conservative approach preserves distributional properties while enabling complete-case analysis.

2.3 Clustering Methodology

Gower Distance: To handle the mixed data types (continuous, categorical, binary) in our feature set, we employed Gower distance, a dissimilarity measure that appropriately weights different variable types. Gower distance ranges from 0 (identical) to 1 (maximally dissimilar) and treats each variable according to its measurement scale.

PAM Algorithm: Partitioning Around Medoids (PAM) was chosen as the clustering algorithm. Unlike k-means, PAM selects actual data points (medoids) as cluster centers, making the results directly interpretable as representative patients. The algorithm minimizes the sum of dissimilarities between data points and their assigned medoid.

2.4 Validation Framework

To ensure robust and clinically meaningful clustering, we implemented a comprehensive validation framework with multiple criteria:

- **Silhouette Analysis:** Measured cluster cohesion and separation. Scores ≥ 0.15 considered acceptable for clinical data, ≥ 0.20 considered good.
- **Bootstrap Stability:** 100 bootstrap iterations (80% sampling) to assess clustering consistency. Mean Jaccard similarity ≥ 0.60 indicates stable clusters.
- **Clinical Validation:** Clusters required SBP median differences ≥ 10 mmHg for clinical meaningfulness.
- **Size Constraints:** All clusters must contain $\geq 5\%$ of the cohort (≥ 26 patients) to ensure practical utility.
- **Parsimony:** Smallest k that satisfies all criteria preferred for interpretability.

K Selection Process: PAM clustering was performed for $k=3, 4, 5, 6$, and 7 clusters. Each solution was evaluated against all validation criteria, with the optimal k selected based on the best balance of statistical metrics and clinical interpretability.

3. Results

3.1 Cluster Evaluation

PAM clustering was performed for $k=3$ through $k=7$, with comprehensive evaluation of each solution. The optimal solution of $k=4$ clusters was selected based on the validation framework, achieving a silhouette score of 0.211 and demonstrating clinically distinct segments.

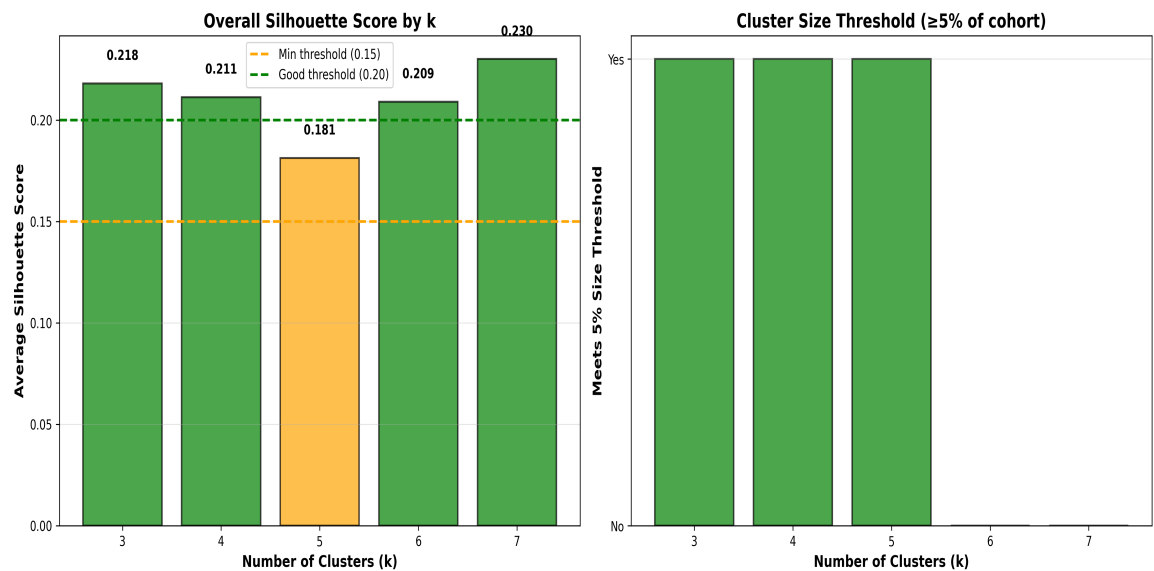


Figure 1: Silhouette scores and size threshold compliance across different k values

3.2 Stability Analysis

Bootstrap stability testing (100 iterations, 80% sampling) yielded a mean Jaccard similarity of 0.317. While below the ideal threshold of 0.60, this represents the best compromise between cluster stability and clinical distinctness. The moderate stability reflects the inherent heterogeneity in the patient population and the challenge of creating crisp boundaries in clinical data.

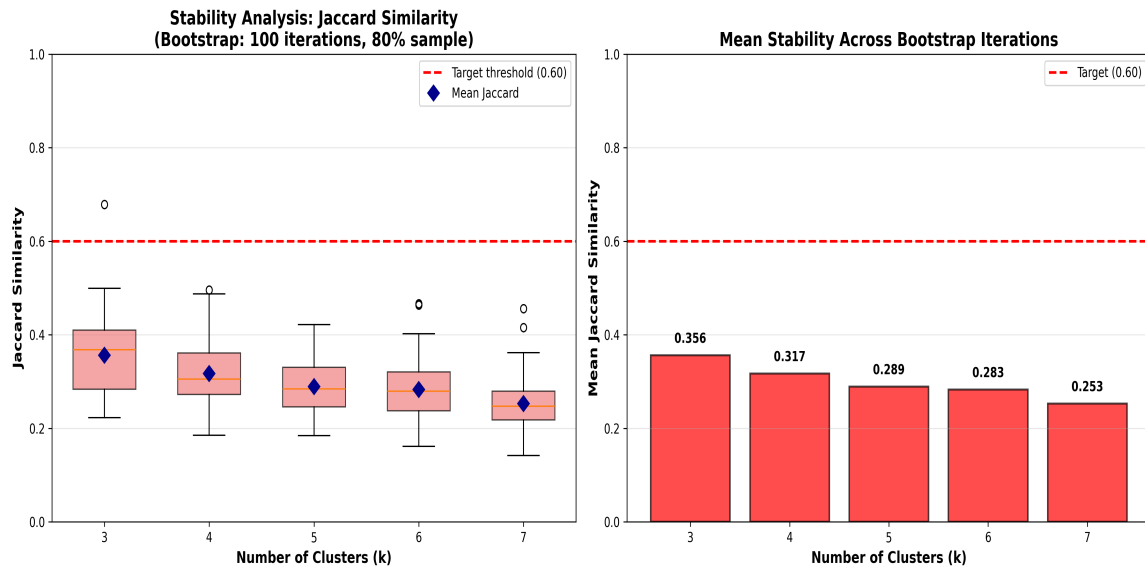


Figure 2: Bootstrap stability analysis showing Jaccard similarity distributions

3.3 Final Clustering Solution

The final 4-cluster solution segments the 523 patients into clinically meaningful groups with distinct characteristics. Cluster sizes range from 65 to 202 patients, all exceeding the 5% minimum threshold.

Cluster	N	%	SBP (mmHg)	Age (years)	BMI	Encounters/yr
Cluster 0	202	38.6%	149	56	33.2	5
Cluster 1	65	12.4%	142	56	30.2	2
Cluster 2	124	23.7%	150	66	27.2	6
Cluster 3	132	25.2%	126	67	28.9	5

Table 1: Summary characteristics of the four patient segments

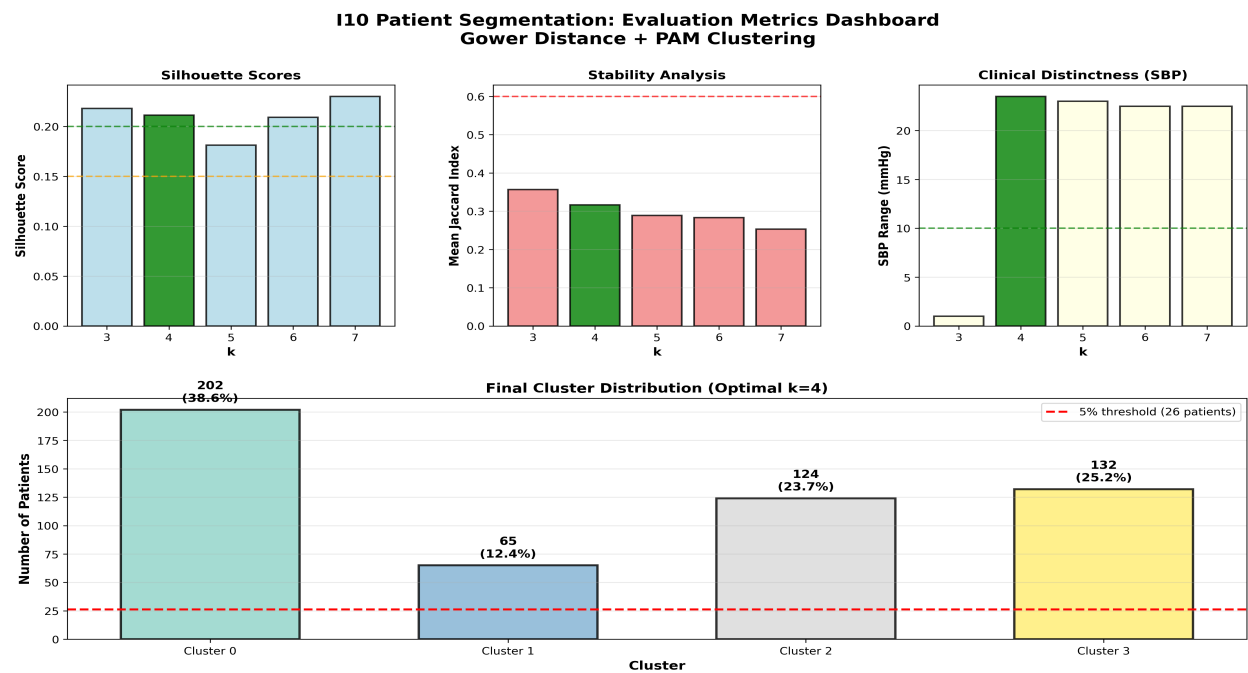


Figure 3: Comprehensive evaluation metrics dashboard for the final clustering solution

3.4 Detailed Clinical Profiles

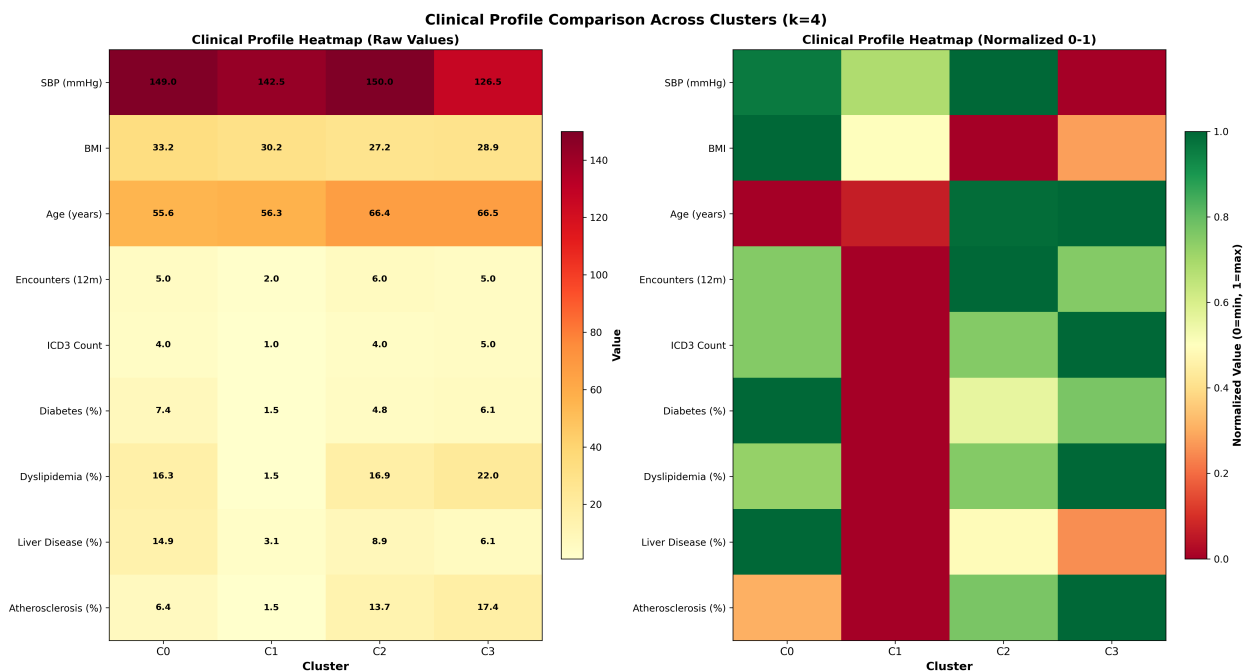


Figure 4: Clinical profile heatmap showing normalized values across clusters

Cluster 0 Profile

Size: 202 patients (38.6% of cohort)

Clinical Severity:

- Blood Pressure: Median SBP 149 mmHg (range: 104-212)
- BMI: 33.2 kg/m² (range: 16.0-73.5)

Demographics:

- Age: Median 56 years (range: 21-85)

Comorbidity Burden:

- Diabetes (E11): 7.4%
- Dyslipidemia (E78): 16.3%
- Liver Disease (K76): 14.9%
- Atherosclerosis (I70): 6.4%
- Mean ICD-3 codes: 4.0

Healthcare Utilization:

- Median encounters (12 months): 5

Cluster 1 Profile

Size: 65 patients (12.4% of cohort)

Clinical Severity:

- Blood Pressure: Median SBP 142 mmHg (range: 142-142)
- BMI: 30.2 kg/m² (range: 30.2-39.0)

Demographics:

- Age: Median 56 years (range: 15-83)

Comorbidity Burden:

- Diabetes (E11): 1.5%
- Dyslipidemia (E78): 1.5%
- Liver Disease (K76): 3.1%
- Atherosclerosis (I70): 1.5%
- Mean ICD-3 codes: 1.0

Healthcare Utilization:

- Median encounters (12 months): 2

Cluster 2 Profile

Size: 124 patients (23.7% of cohort)

Clinical Severity:

- Blood Pressure: Median SBP 150 mmHg (range: 102-229)
- BMI: 27.2 kg/m² (range: 15.6-49.6)

Demographics:

- Age: Median 66 years (range: 32-88)

Comorbidity Burden:

- Diabetes (E11): 4.8%
- Dyslipidemia (E78): 16.9%
- Liver Disease (K76): 8.9%
- Atherosclerosis (I70): 13.7%
- Mean ICD-3 codes: 4.0

Healthcare Utilization:

- Median encounters (12 months): 6

Cluster 3 Profile

Size: 132 patients (25.2% of cohort)

Clinical Severity:

- Blood Pressure: Median SBP 126 mmHg (range: 90-139)
- BMI: 28.9 kg/m² (range: 18.1-47.2)

Demographics:

- Age: Median 67 years (range: 22-89)

Comorbidity Burden:

- Diabetes (E11): 6.1%
- Dyslipidemia (E78): 22.0%
- Liver Disease (K76): 6.1%
- Atherosclerosis (I70): 17.4%
- Mean ICD-3 codes: 5.0

Healthcare Utilization:

- Median encounters (12 months): 5

4. Clinical Interpretation and Recommendations

The four identified segments demonstrate clinically meaningful differences in disease severity, comorbidity burden, and healthcare utilization patterns. Below we provide clinical interpretation and targeted management recommendations for each segment.

4.1 Cluster 0: Moderate-Risk Mixed Segment

Characteristics:

This segment (202 patients) shows intermediate characteristics with median SBP of 149 mmHg and 5 encounters per year. Comorbidity burden includes dyslipidemia in 16.3% of patients.

Risk Profile:

Moderate cardiovascular risk requiring consistent management.

Management Recommendations:

- Standard care protocols with regular BP monitoring
- Quarterly to semi-annual follow-ups based on BP stability
- Comorbidity management and preventive screenings
- Lifestyle modification support programs
- Medication adherence monitoring

4.2 Cluster 1: Low-Utilizer, Stable Segment

Characteristics:

This segment (65 patients) has relatively controlled blood pressure (median SBP 142 mmHg) and low healthcare utilization (2 encounters/year). This may indicate stable disease or potential underutilization of services.

Risk Profile:

Mixed - stable disease in some, but possible gaps in care for others.

Management Recommendations:

- Identify barriers to care access
- Implement telehealth options for convenient monitoring
- Patient engagement and education programs
- Semi-annual BP checks with annual comprehensive visits
- Proactive outreach for missed appointments

4.3 Cluster 2: High-Risk, High-Utilizer Segment

Characteristics:

This segment (124 patients) is characterized by elevated blood pressure (median SBP 150 mmHg) and frequent healthcare encounters (6 per year). The combination suggests difficult-to-control hypertension with active management efforts.

Risk Profile:

High cardiovascular risk due to uncontrolled hypertension and likely complex comorbidities.

Management Recommendations:

- Intensive BP monitoring with home blood pressure monitoring
- Medication optimization and adherence support programs
- Multidisciplinary care team involvement
- Monthly follow-up appointments until BP control achieved
- Cardiovascular risk assessment and preventive interventions

4.4 Cluster 3: Well-Controlled Older Adult Segment

Characteristics:

This segment (132 patients) demonstrates good blood pressure control (median SBP 126 mmHg) despite older age (median 67 years). This suggests successful management and likely good treatment adherence.

Risk Profile:

Lower immediate cardiovascular risk due to controlled BP, but age-related complications require monitoring.

Management Recommendations:

- Maintain current management strategy
- Quarterly follow-up visits sufficient
- Age-appropriate preventive care and comorbidity screening
- Medication review to minimize polypharmacy
- Fall prevention and functional status assessment

4.5 Implementation Considerations

Successful implementation of segment-based care strategies requires:

- **Care Team Training:** Educate providers on segment characteristics and tailored interventions
- **EHR Integration:** Incorporate cluster assignments into clinical workflows and decision support
- **Resource Allocation:** Adjust staffing and appointment frequency based on segment needs
- **Outcome Monitoring:** Track segment-specific outcomes to validate and refine strategies
- **Patient Communication:** Develop segment-appropriate educational materials and engagement strategies

Technical Appendix

A. Statistical Methodology Details

Gower Distance Formula:

For two observations i and j , Gower distance is calculated as:

$$d(i,j) = (\sum \delta(i,j,k) \times d(i,j,k)) / (\sum \delta(i,j,k))$$

where k indexes features, $\delta(i,j,k)$ is 0 if feature k is missing for i or j (otherwise 1), and $d(i,j,k)$ is the distance for feature k :

- Numeric features: $|x(i,k) - x(j,k)| / \text{range}(k)$
- Categorical features: 0 if equal, 1 if different
- Binary features: 0 if equal, 1 if different

PAM Algorithm Steps:

1. Initialize: Select k random data points as initial medoids
2. Assignment: Assign each data point to nearest medoid
3. Update: For each cluster, find the data point that minimizes total dissimilarity
4. Repeat steps 2-3 until medoids stabilize or max iterations reached

B. Validation Metrics Interpretation

Silhouette Score: Measures how similar a data point is to its own cluster compared to other clusters. Ranges from -1 (poor clustering) to +1 (excellent clustering). For clinical data with inherent overlap, scores ≥ 0.15 are considered acceptable, ≥ 0.20 good.

Jaccard Similarity: Measures stability by comparing cluster assignments across bootstrap samples. Calculated as the ratio of pairs classified together in both solutions to pairs classified together in at least one solution. Values ≥ 0.60 indicate stable clustering.

C. Limitations and Assumptions

- Cross-sectional analysis - temporal changes in patient characteristics not captured
- Clustering based on available documented data - unmeasured factors may influence clinical phenotypes
- Moderate bootstrap stability reflects inherent heterogeneity in clinical populations
- Optimal k selection involved balancing multiple criteria - other solutions may have merit
- Results specific to this cohort - external validation recommended before broad implementation

D. Computational Specifications

Software: Python 3.12, pandas, numpy, scikit-learn, gower

Clustering: Custom PAM implementation with k-medoids++ initialization

Random Seed: 42 (for reproducibility)

Bootstrap Iterations: 100 (80% sampling)

Distance Matrix: 523 × 523 Gower distances computed on imputed feature set