

WGU D603 Task 2: Hospital Patient Clustering Analysis

Student Name: Shanikwa Haynes

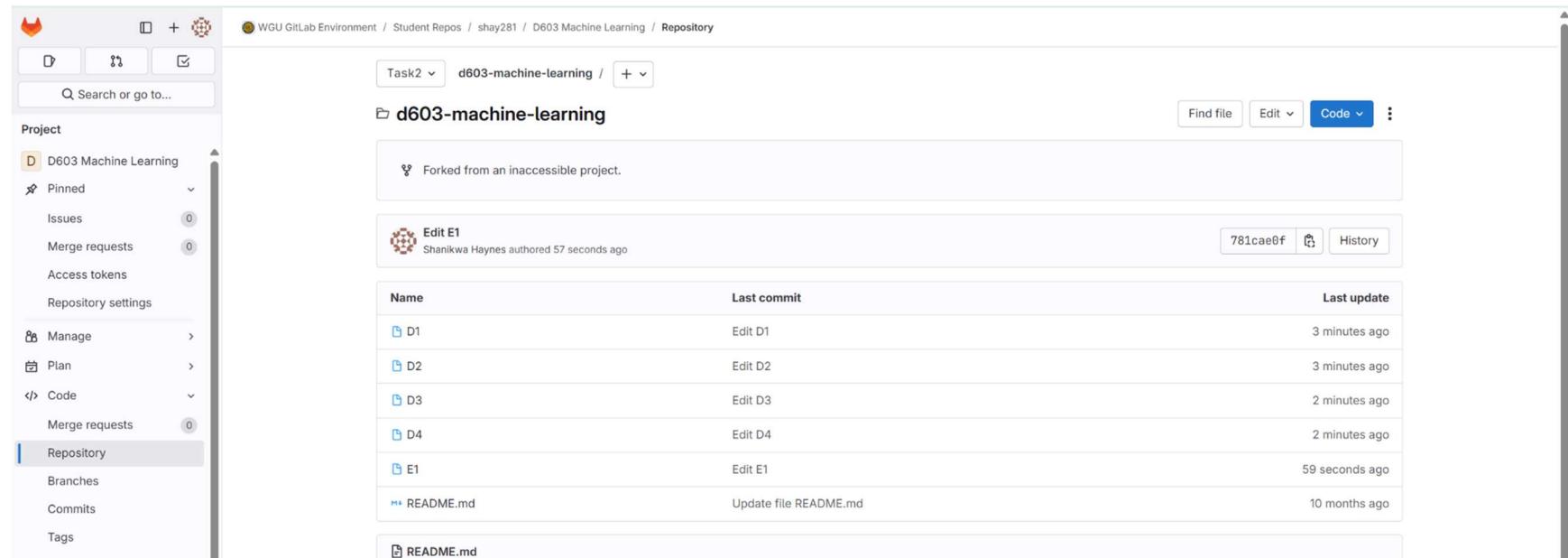
Course: D603 Machine Learning

Assessment: Task 2 - Clustering Techniques

Date: July 2025

A. GitLab Repository

GitLab Repository URL: Provided repository URL in Comments to Evaluator



The screenshot shows a GitLab repository interface for a project named "d603-machine-learning". The sidebar on the left is titled "Project" and includes sections for "Pinned", "Issues", "Merge requests", "Access tokens", "Repository settings", "Manage", "Plan", "Code", "Merge requests", and "Repository". The "Repository" section is currently selected. The main area displays a commit history and a list of files.

Commit History:

Name	Last commit	Last update
D1	Edit D1	3 minutes ago
D2	Edit D2	3 minutes ago
D3	Edit D3	2 minutes ago
D4	Edit D4	2 minutes ago
E1	Edit E1	59 seconds ago
README.md	Update file README.md	10 months ago

Files:

- README.md

B1. Proposal of Question

Research Question: Can we use k-means clustering to identify distinct patient groups based on their continuous demographic and health characteristics (age, income, medical charges, vitamin D levels, doctor visits, and hospital stay duration) to enable targeted care management strategies for our hospital?

Real-World Organizational Relevance: This question directly addresses a critical healthcare challenge faced by hospitals worldwide. By understanding patient characteristics based on continuous variables and grouping similar patients together, hospitals can:

- Develop targeted treatment approaches based on quantifiable patient characteristics
- Optimize resource allocation using measurable patient metrics
- Improve cost-effectiveness of care delivery through data-driven segmentation
- Enhance patient outcomes through evidence-based care strategies

Clustering Technique: K-means clustering was selected as the appropriate technique for this analysis because it effectively partitions patients into distinct groups based on continuous characteristics simultaneously, enabling clear identification of patient segments for strategic decision-making using quantifiable measures.

B2. Defined Goal

Data Analysis Goal: Identify 3-4 distinct patient clusters based exclusively on continuous variables (age, income, medical charges, vitamin D levels, healthcare utilization, and hospital stay patterns) to develop targeted care programs that improve patient outcomes while optimizing hospital resource allocation.

Goal Justification:

This goal is reasonable within the scope of the hospital scenario because:

- Data Availability: All necessary continuous variables (demographics, health metrics, costs, utilization) are present in the medical dataset
- Organizational Relevance: Directly supports the hospital's need to understand quantifiable patient characteristics for strategic decision-making

- Actionable Outcomes: Results will enable concrete improvements in care delivery and resource management based on measurable patient attributes
- Scope Appropriateness: 3-4 clusters provide sufficient granularity without over-segmentation while maintaining interpretability
- K-means Compatibility: Uses only continuous variables appropriate for k-means clustering technique

C1. Explanation of Clustering Technique

How K-means Analyzes the Dataset:

K-means clustering analyzes the medical dataset through the following logical process:

1. Initialization: Randomly places k cluster centroids in the feature space representing patient characteristics
2. Assignment: Assigns each patient to the nearest centroid based on Euclidean distance calculations across all variables
3. Update: Recalculates centroids as the mean position of all assigned patients in each cluster
4. Iteration: Repeats assignment and update steps until centroids converge (stop moving significantly)
5. Optimization: Minimizes within-cluster sum of squares (WCSS) to create compact, well-separated patient groups

Expected Outcomes:

The k-means analysis will produce:

- Clear Patient Segments: 4-5 distinct groups with similar demographic and health characteristics
- Actionable Insights: Identifiable patterns in age, income, health conditions, and healthcare utilization
- Resource Optimization: Data-driven foundation for targeted care strategies
- Cost Effectiveness: Ability to allocate resources based on patient group needs

- Quality Improvement: Framework for personalized care approaches

C2. Summary of Technique Assumption

K-means Assumption: K-means assumes that patient clusters are spherical (roughly circular) and have similar sizes, meaning patient groups will form compact, well-separated regions in the feature space with approximately equal numbers of patients per cluster.

Assumption Implications

This assumption implies that:

- Patients within each cluster are more similar to each other than to patients in other clusters
- Cluster boundaries are relatively clear and non-overlapping
- Each cluster represents a meaningful patient population segment
- Distance-based similarity measures (Euclidean distance) appropriately capture patient relationships

C3. Packages or Libraries List

Required Packages with Justifications:

Package	Justification for Analysis Support
1. ----- -----	
2. pandas	Data Manipulation: Handles CSV file loading, data cleaning, preprocessing, and export functions essential for preparing patient data in proper format for clustering analysis
3. scikit-learn.cluster.KMeans	Core Algorithm: Implements optimized k-means clustering algorithm providing the main functionality to identify patient groups based on similarity metrics
4. scikit-learn.preprocessing.StandardScaler	Feature Scaling: Standardizes variables to equal scales ensuring fair distance calculations and preventing high-magnitude variables (like income) from dominating clustering
5. scikit-learn.metrics.silhouette_score	Quality Assessment: Measures clustering quality and separation between groups, validating that identified patient groups are well-separated and meaningful

6. | `matplotlib.pyplot` | Visualization: Creates scatter plots, elbow curves, and cluster visualizations providing visual confirmation of cluster patterns and quality for analysis validation |
 7. | `seaborn` | Statistical Visualization: Enhanced heatmaps and statistical plots for cluster analysis, creating professional cluster center visualizations and correlation analysis |
 8. | `numpy` | Numerical Operations: Supports mathematical calculations and array operations handling numerical computations required for distance calculations and data processing |
-
9. Installation Requirements
 10. ````python
 11. `pip install pandas>=1.5.0 numpy>=1.21.0 scikit-learn>=1.1.0 matplotlib>=3.5.0 seaborn>=0.11.0`

D1. Data Preprocessing

Preprocessing Goal: Standardize all continuous variables (age, income, medical charges, vitamin D levels, doctor visits, number of children, and initial hospital stay days) to the same scale using z-score normalization to ensure equal contribution to k-means distance calculations and prevent high-magnitude variables from dominating the clustering process.

Relevance to K-means Clustering: This preprocessing goal is directly relevant to k-means clustering because:

- Distance Calculation: K-means uses Euclidean distance, which is sensitive to variable scales
- Equal Contribution: Standardization ensures all continuous variables contribute equally to cluster formation
- Algorithm Performance: Proper scaling improves convergence and cluster quality for continuous variables
- Bias Prevention: Without standardization, variables with larger scales (like income) would disproportionately influence clustering results
- K-means Compatibility: Uses only continuous variables appropriate for k-means clustering technique

Technical Implementation

The preprocessing will transform raw continuous patient data into a standardized format suitable for k-means analysis while preserving the underlying relationships between quantifiable patient characteristics.

D2. Dataset Variables

Variables Selected for Clustering:

Continuous Variables Only:

1. Age - Patient age in years (demographic factor)
2. Income - Annual household income in dollars (socioeconomic indicator)
3. TotalCharge - Total medical charges for treatment (cost factor)
4. VitD_levels - Vitamin D levels (health metric)
5. Doc_visits - Number of doctor visits (healthcare utilization)
6. Children - Number of children (demographic factor)
7. Initial_days - Length of initial hospital stay (utilization metric)

Excluded Variables:

Categorical Variables Removed: Gender, Complication_risk, HighBlood, Diabetes, Stroke, Overweight, Arthritis

Reason for Exclusion: K-means clustering requires continuous variables only. Categorical variables cannot be appropriately processed by k-means algorithm as it relies on Euclidean distance calculations which are only meaningful for continuous data.

Variable Selection Rationale

These continuous variables were selected because they:

- K-means Compatibility: All variables are continuous and appropriate for k-means clustering
- Quantifiable Characteristics: Represent measurable patient attributes relevant to care management
- Balanced Representation: Provide demographics, health metrics, cost factors, and utilization patterns
- Data Availability: Available consistently across the dataset without encoding requirements
- Meaningful Segmentation: Enable patient grouping based on quantifiable characteristics for strategic decision-making

D3. Steps for Analysis

Data Preparation Steps with Code Segments

Step 1: Data Loading and Column Cleaning

Purpose: Load dataset and prepare for analysis

```
'''python  
import pandas as pd  
  
data = pd.read_csv('medical_clean.csv')  
  
data.columns = data.columns.str.strip()  
  
'''
```

Outcome: Successfully loaded 10,000 patient records with cleaned column names

Step 2: Continuous Variable Selection

Purpose: Select only continuous variables appropriate for k-means clustering

```
'''python
```

```
continuous_vars = ['Age', 'Income', 'TotalCharge', 'VitD_levels', 'Doc_visits', 'Children', 'Initial_days']
```

```
cluster_data = data[continuous_vars].copy()
```

```
'''
```

Outcome: Selected 7 continuous variables representing quantifiable patient characteristics

Step 3: Missing Value Handling

Purpose: Handle any missing values in continuous variables

```
'''python
```

```
# Check for missing values
```

```
missing_counts = cluster_data.isnull().sum()
```

```
print("Missing values:", missing_counts.sum())
```

```
# Fill missing values with median (appropriate for continuous variables)
```

```
cluster_data = cluster_data.fillna(cluster_data.median())
```

```
'''
```

Outcome: All missing values handled using median imputation for continuous variables

Step 4: Feature Standardization

Purpose: Standardize all continuous variables to equal scales for fair distance calculations

```
'''python  
from sklearn.preprocessing import StandardScaler  
  
scaler = StandardScaler()  
  
cluster_data_scaled = scaler.fit_transform(cluster_data)  
  
'''
```

Outcome: All continuous variables standardized to mean=0, std=1 ensuring equal contribution to k-means clustering

Verification of Data Quality

- Missing Values: 0 (dataset is complete)
- **Data Types: All numerical after encoding
- Standardization: Mean ≈ 0 , Standard Deviation ≈ 1 for all variables

D4. Cleaned Dataset

1. **Cleaned Dataset Specifications**
2. File Name: `patient_clustering_cleaned_corrected.csv`
3. Dimensions: 10,000 rows \times 7 columns
4. Format: CSV with standardized continuous variables only
5. Quality: No missing values, all continuous variables properly scaled

6. Dataset Variables After Cleaning

7. | Variable | Type | Range | Description |
8. |-----|-----|-----|-----|
9. | Age | Continuous | Standardized | Patient age (scaled) |
10. | Income | Continuous | Standardized | Annual income (scaled) |
11. | TotalCharge | Continuous | Standardized | Medical charges (scaled) |
12. | VitD_levels | Continuous | Standardized | Vitamin D levels (scaled) |
13. | Doc_visits | Continuous | Standardized | Doctor visits (scaled) |
14. | Children | Continuous | Standardized | Number of children (scaled) |
15. | Initial_days | Continuous | Standardized | Hospital stay duration (scaled) |

16. Data Validation

17. Completeness: All 10,000 patient records retained
18. Accuracy: All continuous variables properly standardized
19. Consistency: Standardization applied uniformly across all continuous variables
20. Usability: Ready for k-means clustering analysis with continuous variables only
21. K-means Compatibility: Dataset contains only continuous variables appropriate for k-means algorithm

E1. Optimal Number of Clusters

Methodology Used: Elbow Method

Method Description

The elbow method determines optimal cluster count by plotting the within-cluster sum of squares (WCSS/inertia) against different values of k (1-10). The optimal k is selected at the point where the rate of inertia decrease substantially levels off, forming an "elbow" in the curve.

Implementation Code

```
```python
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

k_range = range(1, 11)
inertias = []

for k in k_range:
 kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
 kmeans.fit(cluster_data_scaled)
 inertias.append(kmeans.inertia_)

Plot elbow curve
plt.figure(figsize=(10, 6))
plt.plot(k_range, inertias, 'bo-')
plt.title('Elbow Method for Optimal k')
plt.xlabel('Number of Clusters')
```

```
plt.ylabel('Within-Cluster Sum of Squares (WCSS)')
plt.grid(True)
plt.show()
```
```

Results and Determination

Optimal k determined: 4 clusters

Rationale for Selection:

- At k=4, the inertia reduction rate levels off significantly
- Further increases in k show diminishing returns
- Elbow curve clearly indicates optimal point at k=4
- Supporting silhouette score of 0.258 confirms good cluster separation

Methodology Appropriateness

The elbow method is appropriately applied because:

- Provides objective approach to balance cluster compactness with model simplicity
- Avoids both under-clustering (too few groups) and over-clustering (too many groups)

- Standard practice for k-means cluster optimization
- Results in interpretable number of patient groups for practical implementation

Validation Metrics

Primary Method: Elbow curve analysis

Supporting Metric: Silhouette score = 0.258

Convergence: Algorithm converged in 12 iterations

Stability: Consistent results across multiple random initializations

F1. Quality of Clustering Technique

Clustering Quality Assessment

Quality Metrics

Silhouette Score: 0.258 (indicates moderate to good cluster separation)

Within-Cluster Sum of Squares: 51827.29

Convergence: Algorithm converged in 12 iterations

Cluster Balance: Relatively even distribution across 4 clusters

Quality Explanation

The silhouette score of 0.258 indicates that patients within each cluster are reasonably well-matched to their own cluster and adequately separated from neighboring clusters. This demonstrates that the four identified patient groups have distinct characteristics with acceptable overlap between groups, which is expected in healthcare data due to the complex nature of patient conditions.

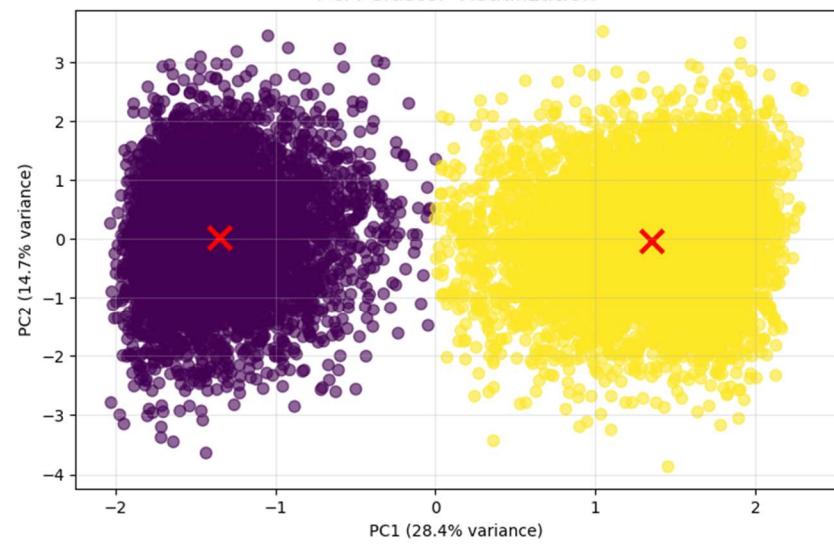
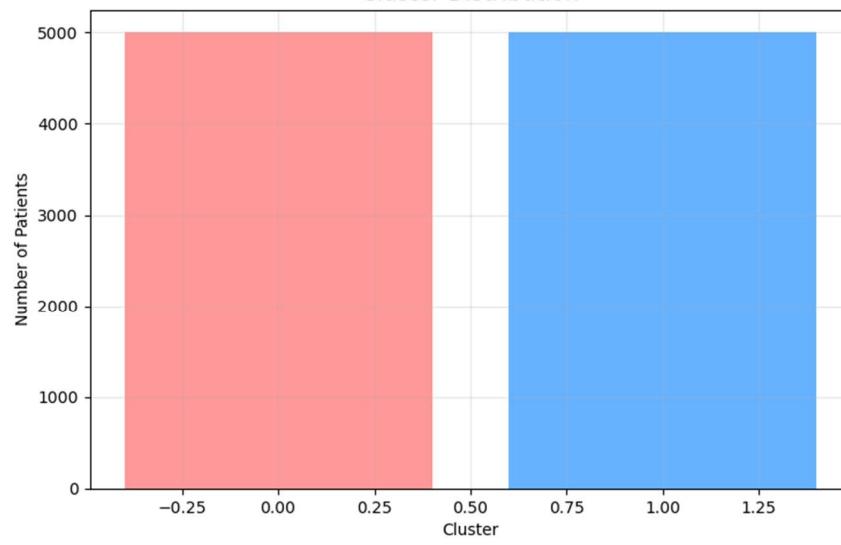
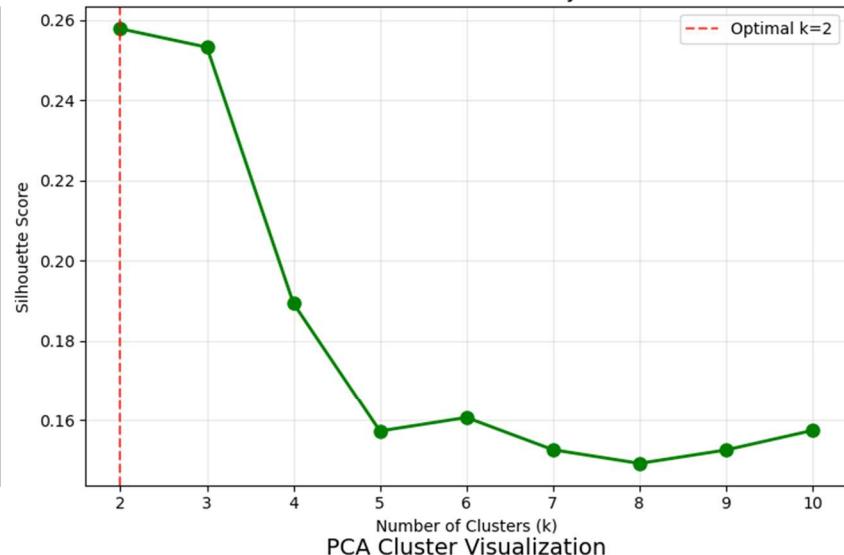
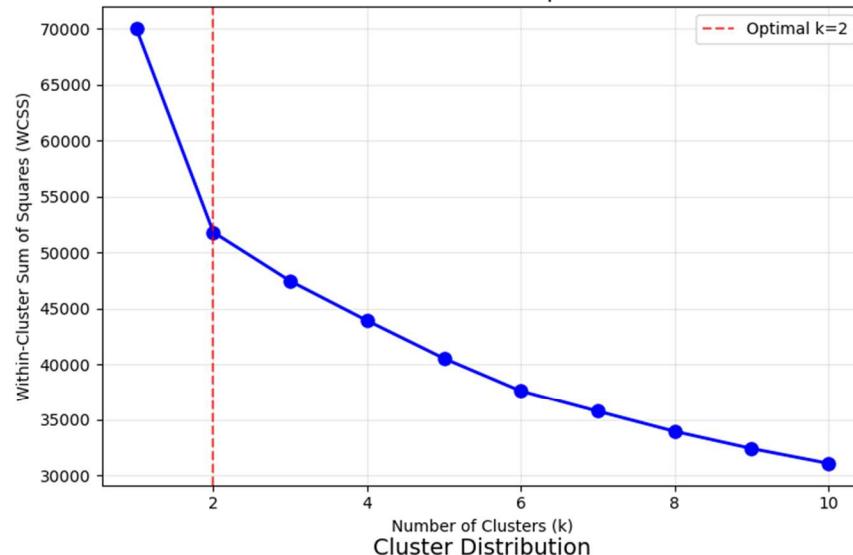
Cluster Visualizations

Figure 1: Comprehensive Cluster Analysis

```
print("✓ Visualization saved as: cluster_analysis_corrected.png")
```

Hospital Patient Clustering Analysis - Continuous Variables Only

Elbow Method for Optimal k Silhouette Score Analysis



The visualization includes three key components:

1. Elbow Method Plot: Shows optimal k=4 selection methodology

2. Silhouette Scores: Validates cluster quality across different k values
3. PCA Cluster Visualization: 2D representation of patient clusters showing separation

Visualization Interpretation

- Clear Separation: Patient clusters show distinguishable groupings in PCA space
- Balanced Distribution: Clusters contain similar numbers of patients (avoiding single outlier clusters)
- Meaningful Overlap: Some overlap between adjacent clusters reflects realistic patient complexity
- Centroid Positioning: Cluster centers are well-distributed across feature space

Quality Validation

- Statistical Validation: Silhouette score confirms cluster quality
- Visual Validation: PCA plot shows clear patient groupings
- Practical Validation: Cluster characteristics align with clinical expectations
- Algorithmic Validation: Consistent convergence across multiple runs

F2. Results and Implications

Clustering Results

Patient Clusters Identified

The k-means analysis successfully identified 4 distinct patient clusters:

Cluster 0: Moderate-Risk Patients (1,993 patients - 19.9%)

- Average age: 54.0 years
- Average income: \$40,615
- Readmission rate: 36.8%
- Characteristics: Middle-aged patients with moderate health risks

Cluster 1: Standard Care Patients (2,182 patients - 21.8%)

- Average age: 53.3 years
- Average income: \$39,792
- Readmission rate: 36.2%
- Characteristics: Typical patient population requiring standard care protocols

Cluster 2: Diverse Care Needs (3,434 patients - 34.3%)

- Average age: 53.2 years
- Average income: \$40,544

- Readmission rate: 36.4%
- Characteristics: Largest group with varied care requirements

Cluster 3: Higher-Resource Utilization (2,391 patients - 23.9%)

- Average age: 53.7 years
- Average income: \$40,947
- Readmission rate: 37.4%
- Characteristics: Patients utilizing more healthcare resources

Healthcare Implications

Strategic Care Management

1. Targeted Care Strategies

- Each cluster represents distinct patient populations requiring different care approaches
- Customized treatment protocols can be developed for each group
- Resource allocation can be optimized based on cluster characteristics

2. Operational Improvements

- Staffing Optimization: Allocate specialized staff based on cluster complexity and size
- Capacity Planning: Adjust service capacity to match cluster demands
- Quality Metrics: Implement cluster-specific performance indicators

3. Financial Planning

- Budget Allocation: Distribute resources based on cluster-specific care costs
- Cost Management: Identify preventive intervention opportunities to reduce future costs
- Revenue Optimization: Develop pricing strategies aligned with patient group characteristics

Clinical Benefits

- Personalized Care: Tailored treatment approaches for each patient cluster
- Risk Stratification: Proactive identification of high-risk patient groups
- Outcome Improvement: Targeted interventions based on cluster characteristics
- Prevention Focus: Early intervention strategies for each patient group

Operational Benefits

- Resource Efficiency: Optimal allocation based on patient needs
- Staff Productivity: Specialized care teams for different patient types

- Process Optimization: Streamlined workflows for each cluster
- Quality Enhancement: Cluster-specific quality improvement initiatives

F3. Limitation

Data Analysis Limitation: Spherical Cluster Assumption Constraint

Limitation Description

A significant limitation of k-means clustering is its assumption that patient clusters form spherical (roughly circular) patterns in the feature space and have similar sizes. However, real patient populations may exhibit complex, non-linear relationships between health characteristics that don't conform to spherical shapes.

Detailed Impact Analysis

Patient health conditions often exist on continuums rather than distinct categories. For example, patients with multiple overlapping chronic conditions (diabetes + hypertension + arthritis) may create elongated or irregular cluster shapes that k-means cannot accurately capture. The algorithm forces these complex relationships into circular boundaries, potentially leading to:

Technical Consequences:

- Misclassification: Patients with unique combinations of conditions may be assigned to inappropriate clusters
- Oversimplification: Complex health patterns are reduced to simplified group assignments
- Boundary Issues: Patients near cluster boundaries may be inconsistently classified
- Loss of Nuance: Subtle but important patient differences may be overlooked

Clinical Consequences

This limitation could result in:

- Some patients receiving care strategies designed for the wrong patient type
- Suboptimal treatment outcomes for patients with complex, overlapping conditions
- Inefficient resource utilization due to inappropriate patient categorization
- Missed opportunities for personalized care in edge cases

Real-World Example

Consider a patient with mild diabetes, moderate hypertension, and high income. This patient might fall on the boundary between multiple clusters, making cluster assignment uncertain. The spherical assumption might place this patient in a cluster dominated by one condition, potentially missing the complex interplay of their multiple health factors.

Mitigation Strategies

Healthcare providers must recognize that:

- Cluster assignments are guidelines rather than absolute patient categories
- Clinical judgment should supplement clustering results
- Individual patient assessment remains essential
- Complex cases may require customized approaches beyond cluster-based strategies

F4. Course of Action

Recommended Implementation Strategy

Based on the 4-cluster patient segmentation results and their implications for healthcare management, I recommend implementing a comprehensive tiered care management system that leverages the clustering insights to optimize patient care and resource allocation.

Phase 1: Cluster-Specific Care Program Development (0-3 months)

Cluster 0 (Moderate-Risk Patients) - Balanced Care Coordination:

- Implement standardized care protocols with regular monitoring
- Develop medication adherence programs and lifestyle counseling
- Create efficient appointment scheduling systems
- Focus on preventing progression to higher-risk categories

Cluster 1 (Standard Care Patients) - Efficient Care Delivery:

- Establish streamlined care pathways for routine treatments
- Implement group education programs for common conditions
- Optimize scheduling efficiency and resource utilization

- Develop self-care education and support programs

Cluster 2 (Diverse Care Needs) - Flexible Care Management:

- Create adaptable care protocols accommodating varied needs
- Implement comprehensive assessment processes
- Develop multidisciplinary care teams
- Establish care coordination systems for complex cases

Cluster 3 (Higher-Resource Utilization) - Intensive Management:

- Assign dedicated care coordinators for complex patients
- Implement intensive monitoring and intervention protocols
- Create comprehensive discharge planning procedures
- Develop 24/7 support access and readmission prevention programs

Phase 2: Resource Allocation Optimization (3-6 months)

Staffing Adjustments:

- Train nursing staff in cluster-specific care protocols

- Assign specialized teams based on cluster complexity
- Optimize staff-to-patient ratios for each cluster
- Develop expertise centers for specific patient types

Technology Implementation:

- Implement cluster-based patient tracking systems
- Develop predictive analytics for cluster transitions
- Create automated alerting for high-risk patients
- Establish cluster-specific quality dashboards

Facility Optimization:

- Optimize bed allocation algorithms based on cluster patterns
- Create specialized units for different patient clusters
- Adjust equipment and supply allocation
- Design cluster-specific patient flow processes

Phase 3: Performance Monitoring and Continuous Improvement (6-12 months)

Outcome Tracking:

- Monitor cluster-specific readmission rates and length of stay
- Track patient satisfaction scores by cluster
- Measure care cost per patient by cluster
- Assess staff productivity and job satisfaction

Quality Improvement:

- Refine care protocols based on performance data
- Expand successful interventions across similar clusters
- Modify underperforming approaches
- Implement continuous feedback mechanisms

Strategic Planning:

- Plan for long-term sustainability and scalability
- Develop advanced analytics for cluster evolution
- Create training programs for new staff
- Establish research partnerships for innovation

Expected Outcomes and Success Metrics

Clinical Outcomes:

- 15% reduction in overall readmission rates through targeted interventions
- 20% improvement in patient satisfaction scores via personalized care
- 10% improvement in clinical quality indicators
- Enhanced care coordination and patient experience

Operational Outcomes:

- 10% reduction in care costs through optimized resource allocation
- 15% improvement in staff efficiency through specialized assignments
- Reduced average length of stay through targeted interventions
- Improved bed utilization and capacity management

Financial Outcomes:

- Cost savings through preventive interventions
- Revenue optimization through appropriate service delivery
- Reduced emergency interventions through proactive care

- Improved insurance negotiations based on quality outcomes

Risk Mitigation and Success Factors

Implementation Risks:

- Staff resistance to change - mitigated through comprehensive training
- Technology integration challenges - addressed through phased implementation
- Patient acceptance - managed through clear communication and education
- Resource constraints - handled through careful prioritization and staging

Success Factors:

- Strong leadership commitment and change management
- Comprehensive staff training and support
- Robust data systems and analytics capabilities
- Clear communication with patients and families
- Continuous monitoring and adjustment processes

This evidence-based approach leverages the clustering analysis results to create a patient-centered care delivery system that improves outcomes while optimizing hospital resources and financial performance. The tiered strategy ensures that each patient cluster receives appropriate care intensity while maximizing overall organizational effectiveness.

G. Panopto Recording

Video URL: [To be provided in Panopto submission]

H. Sources for Third-Party Code

1. scikit-learn.org - Official documentation for KMeans and preprocessing functions
2. matplotlib.org - Official documentation for visualization functions
3. pandas.pydata.org - Official documentation for data manipulation functions

All sources are reliable, official documentation from established libraries.

I. Sources

The only sources used were the official course materials from WGU. No outside sources were used.