

Using K-means Clustering to Identify High-Risk Heart Patients Based on Blood Pressure and Cholesterol Data

Executive Summary:

The aim of this MATLAB project is to use k-means clustering to group 20 patients based on their blood pressure and cholesterol levels, into those with high risk of heart attack and those with low risk of heart attack. The project scope includes problem definition, technical review, design requirements, design description, evaluation, references, and appendices.

The problem is approached by clustering the patients based on their blood pressure and cholesterol levels using the k-means algorithm. The technical review considers the theory behind k-means clustering and its implementation in MATLAB. The design requirements outline the necessary input data, the desired output, and the constraints on the algorithm.

The design description provides an overview of the approach, including the preprocessing of the data, implementation of the k-means algorithm, and visualization of the results. The detailed description outlines the implementation of each step in the algorithm and the MATLAB code used to achieve the results. The use section provides instructions on how to run the code and interpret the results.

The evaluation section provides an overview of the prototype, testing methodology, and results of the clustering. The results are assessed based on their accuracy in correctly classifying patients into high and low-risk groups. Finally, the appendices include additional details on the code implementation and any supporting materials.

Overall, the project aims to demonstrate the effectiveness of k-means clustering for grouping patients based on their blood pressure and cholesterol levels, and to provide a practical implementation of the algorithm in MATLAB.

Table of Contents:

1. Problem Definition
 - 1.1 Problem Scope
 - 1.2 Technical Review
 - 1.3 Design Requirements
2. Design Description
 - 2.1 Overview
 - 2.2 Detailed Description
 - 2.3 Use
3. Evaluation
 - 3.1 Overview
 - 3.2 Prototype
 - 3.3 Testing and Results
 - 3.4 Assessment
4. References
5. Appendices

1. Problem Definition:

1.1 Problem Scope:

Heart disease is one of the leading causes of death globally, and high blood pressure and cholesterol levels are among the significant risk factors for heart disease. Therefore, it is crucial to identify individuals with high risk of heart attack and provide them with appropriate care and treatment. The project aims to use k-means clustering to group patients based on their blood pressure and cholesterol levels to identify those with high risk of heart attack and those with low risk. The project seeks to provide a solution that can assist medical professionals in identifying patients who need preventive care and treatment for heart disease. The project will utilize MATLAB to implement the clustering algorithm, and the results will be presented in a clear and understandable format.

1.2 Technical Review:

The k-means clustering algorithm is a widely used unsupervised learning method in machine learning and data mining. It is used to partition data into k groups based on the similarity between the data points. In this project, we will use k-means clustering to group patients into high-risk and low-risk categories based on their blood pressure and cholesterol levels.

1.3 Design Requirements:

The project requires the following design requirements:

- Import the data about blood pressure and cholesterol for 20 patients
- Preprocess the data to remove any missing values or outliers
- Use k-means clustering to group the patients into two clusters: high-risk and low-risk
- Visualize the results of the clustering to see which patients belong to each cluster
- Evaluate the effectiveness of the clustering algorithm

2. Design Description:

2.1 Overview:

The project will use MATLAB to import the data, preprocess it, and apply k-means clustering to group patients into two clusters based on their blood pressure and cholesterol levels. The clustering results will be visualized using MATLAB plots.

2.2 Detailed Description:

Step 1: Import the data about blood pressure and cholesterol for 20 patients

Step 2: Preprocess the data to remove any missing values or outliers

Step 3: Apply k-means clustering to the preprocessed data with $k = 2$ to group the patients into two clusters: high-risk and low-risk

Step 4: Visualize the clustering results using MATLAB plots to see which patients belong to each cluster

Step 5: Evaluate the effectiveness of the clustering algorithm by calculating metrics such as the silhouette score and within-cluster sum of squares.

2.3 Use:

The k-means clustering algorithm is a powerful tool in identifying patterns and grouping similar data points together. In the context of healthcare, it can be used to identify patients who are at a high risk of heart attack based on their blood pressure and cholesterol levels. By using this algorithm, healthcare professionals can more easily identify and monitor patients who may require more intensive treatment or lifestyle changes.

This project has the potential to improve patient outcomes and reduce healthcare costs. By identifying patients at a high risk of heart attack earlier, healthcare professionals can intervene with preventative measures such as medication and lifestyle changes, potentially reducing the incidence and severity of heart attacks. This can lead to improved patient health outcomes and a reduction in the associated healthcare costs.

3. Evaluation:

3.1 Overview:

The project's evaluation involves testing the clustering algorithm's effectiveness and the results' accuracy. The project will calculate metrics such as the silhouette score and within-cluster sum of squares to evaluate the clustering algorithm's effectiveness.

3.2 Prototype:

A prototype of the project will be developed to import and preprocess the data, apply k-means clustering, and visualize the results. The prototype will also calculate metrics such as the silhouette score and within-cluster sum of squares to evaluate the clustering algorithm's effectiveness.

3.3 Testing and Results:

The project will be tested using sample data to determine if it accurately groups patients into high-risk and low-risk categories based on their blood pressure and cholesterol levels. The results will be evaluated based on the accuracy of the clustering algorithm and the calculated metrics.

4. References:

1. Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4). springer.
2. MacQueen, J. (1967). Some methods for classification and analysis of multivariate observations. In Proceedings of the fifth Berkeley symposium on mathematical statistics and probability (Vol. 1, No. 14, pp. 281-297).
3. Matlab Documentation: Clustering (2023). Retrieved from <https://www.mathworks.com/help/stats/clustering.html>

5. Appendices:

1. MATLAB Code used for the project
2. Sample output of clustering results

APPENDIX 1: MATLAB Code used for the project

```
% Load data from file
data = [120 180; 130 190; 140 210; 120 170; 110 160; 130 200; 140 190; 130 180;
120 190; 130 210; 140 200; 120 180; 110 170; 130 190; 140 210; 130 170; 120 200;
130 190; 140 200; 120 180];

% Set k value for k-means clustering
k = 2;

% Run k-means clustering algorithm
[idx, centroids] = kmeans(data, k);

% Print the clusters and their corresponding patients
for i = 1:k
    cluster_idx = find(idx == i);
    fprintf('Cluster %d:\n', i);
    fprintf('Patients: %s\n', num2str(cluster_idx));
    fprintf('Centroid: %s\n\n', num2str(centroids(i,:)));
end

% Plot the clusters and their corresponding patients
figure;
scatter(data(idx==1,1), data(idx==1,2), 'r');
hold on;
scatter(data(idx==2,1), data(idx==2,2), 'b');
scatter(centroids(:,1), centroids(:,2), 'k', 'filled');
legend('Low Risk', 'High Risk', 'Centroids');
xlabel('Blood Pressure');
ylabel('Cholesterol');
title('Blood Pressure and Cholesterol Clustering');
```

APPENDIX 2: Sample output of clustering results

Code Output:

```
Cluster 1:
Patients: 3 6 7 10 11 15 17 19
Centroid: 135      202.5

Cluster 2:
Patients: 1 2 4 5 8 9 12 13 14 16 18 20
Centroid: 122.5    179.1667
```

Plot:

