

Multivariate Analysis of Emergency Department Patient Satisfaction

Claudio Mazzi

Department of Computer Science, University of Pisa
MeS Laboratory, Sant'Anna School for Advanced Studies, Pisa
`claudio.mazzi@santannapisa.it`

Applied Statistical Modelling

A.Y. 2025–2026

1 Introduction

This document presents a complete multivariate analysis workflow applied to a real-world patient satisfaction survey collected in Emergency Departments (EDs). The dataset (`DATA_PS.dta`) contains responses from 5,919 patients to a structured questionnaire covering different domains of the ED experience, such as waiting times, comfort, communication, and evaluations of medical and nursing staff.

The overall objective of this toolkit is didactic: it shows how to combine *factor analysis* and *cluster analysis* to identify latent dimensions of satisfaction and empirically derive patient experience profiles. The main steps are:

- preprocessing and selection of the questionnaire items of interest;
- separate factor analyses for conceptually coherent blocks of items;
- extraction of domain-specific latent factor scores (waiting time, comfort, etc.);
- hierarchical clustering (Ward's method) on the factor scores;
- validation via k-means clustering and visualisation through PCA.

2 Data Description

The `DATA_PS.dta` file contains ED patient satisfaction data obtained from a structured survey. Each row corresponds to one patient; columns include administrative variables and a large number of questionnaire items. In particular, we focus on the original item variables with suffix `_orig`, which store the numeric codes of the responses (e.g. Likert-type scales).

For this exercise, we concentrate on the following groups of items:

- **Triage and waiting time:** `d009_orig`, `d010_orig`, `d016_orig`;
- **Comfort and cleanliness:** `d017_orig`, `d018_orig`;
- **Assistance and communication:** `d021_orig`, `d022_orig`, `d023_orig`, `d024_orig`;
- **Doctors:** `d026_orig`, `d027_orig`, `d028_orig`, `d029_orig`;
- **Nurses:** `d031_orig`, `d032_orig`, `d033_orig`, `d034_orig`.

Each group reflects a specific conceptual domain of ED experience and is measured by two to four ordinal items (e.g. 1–4, 1–6 or 1–7 response categories).

3 Methodological Framework

3.1 Factor Analysis

Factor Analysis (FA) is a multivariate technique used to model the covariance structure of a set of observed variables in terms of a smaller number of unobserved (latent) factors. In the common factor model, each observed variable X_j is expressed as:

$$X_j = \mu_j + \lambda_{j1}F_1 + \lambda_{j2}F_2 + \cdots + \lambda_{jm}F_m + \varepsilon_j, \quad (1)$$

where:

- μ_j is the mean of variable X_j ,
- F_1, \dots, F_m are latent common factors,
- λ_{jk} are the factor loadings,
- ε_j is a specific error term unique to X_j .

In this application, we specify *one* common factor per domain (e.g. one factor for triage and waiting, one for comfort), and we are interested primarily in:

1. verifying that the items within each domain share a strong common dimension;
2. computing factor scores for each patient, representing their position along that latent dimension.

We use *principal axis factoring* (R function `fa()` with `fm = "pa"`). Because the items are ordinal and may contain missing values, we first compute a pairwise correlation matrix for each item block and then run FA on that correlation matrix.

3.2 From Items to Latent Domain Scores

For each domain, we:

1. select the corresponding `_orig` items;
2. compute a pairwise-complete Pearson correlation matrix;
3. perform one-factor FA using principal axis factoring;
4. obtain factor scores for each respondent.

The resulting five scores are:

- **attesa**: latent factor for triage and waiting time;
- **comfort**: latent factor for comfort and cleanliness;
- **assistenza**: latent factor for assistance and communication;
- **medici**: latent factor for doctor-related experience;
- **infermieri**: latent factor for nurse-related experience.

These scores are continuous variables and serve as a compact, noise-reduced representation of the original questionnaire items.

3.3 Hierarchical Cluster Analysis (Ward's Method)

Cluster analysis is an unsupervised learning technique whose goal is to partition a set of observations into groups (clusters) that are internally homogeneous and externally heterogeneous. Here we apply hierarchical agglomerative clustering with Ward's method.

In hierarchical agglomerative clustering:

- we start with each observation as its own cluster;
- at each step, we merge the two clusters that are most similar according to a linkage criterion;
- this process continues until all observations are merged into a single cluster.

Ward's method uses a sum-of-squares criterion: at each step, it merges the two clusters that result in the minimum increase in the total within-cluster sum of squares. This tends to produce clusters that are compact and spherical in shape. The results are visualised via a dendrogram, which displays the nested structure of clusters and helps determine a reasonable number of groups (e.g. by inspecting the height at which branches are merged).

3.4 k-means Clustering and PCA for Visualisation

To validate the structure suggested by hierarchical clustering, we also apply *k-means* clustering on the same factor score matrix. k-means is a partitioning algorithm that:

1. initializes k cluster centroids;
2. assigns each observation to the nearest centroid;
3. recomputes centroids as the mean of assigned points;
4. iterates until convergence.

In our case, both the dendrogram and heuristic criteria support a three-cluster solution, which we further explore using k-means with $k = 3$.

For visualisation, the function `fviz_cluster()` internally applies Principal Component Analysis (PCA) to reduce the five-dimensional factor score space into two principal components (PC1 and PC2). The resulting scatterplot shows cluster membership in this reduced space, where PC1 typically captures the main gradient of overall satisfaction across domains.

4 R Implementation

4.1 Loading the Dataset and Preparing Item Variables

Listing 1: Loading the dataset and preparing *_orig items

```
1 library(haven)
2 library(dplyr)
3 library(psych)
4 library(factoextra)
5 library(cluster)
6
7 # Load the original dataset in Stata format
8 ps <- read_dta("DATA\\_PS.dta")
```

```

9
10 # Convert labelled variables into factors
11 ps <- ps %>% mutate(across(where(is.labelled), ~as_factor(.)))
12
13 # Convert *_orig item variables to numeric codes
14 ps_num <- ps %>%
15   mutate(across(matches("^d[0-9]+_orig$"),
16             ~as.numeric(as.character(.))))

```

4.2 Factor Analysis Helper Function and Domain Scores

Listing 2: Factor analysis helper function and factor scores

```

1 # Helper function: 1-factor principal axis FA + factor scores
2 compute_factor <- function(data, vars) {
3   # Pairwise-complete correlation matrix
4   R <- cor(data[, vars], use = "pairwise.complete.obs")
5   # One-factor principal axis FA
6   fa_res <- fa(r = R, nfactors = 1, fm = "pa")
7   # Factor scores for all observations
8   scores <- factor.scores(data[, vars], fa_res)$scores[, 1]
9   return(scores)
10 }
11
12 # TRIAGE & WAITING TIME: d009_orig, d010_orig, d016_orig
13 ps_num$attesa <- compute_factor(ps_num,
14                                c("d009_orig", "d010_orig", "d016_orig"))
15
16 # COMFORT & CLEANLINESS: d017_orig, d018_orig
17 ps_num$comfort <- compute_factor(ps_num,
18                                 c("d017_orig", "d018_orig"))
19
20 # ASSISTANCE & COMMUNICATION: d021_orig, d022_orig, d023_orig, d024_orig
21 ps_num$assistenza <- compute_factor(ps_num,
22                                    c("d021_orig", "d022_orig",
23                                      "d023_orig", "d024_orig"))
24
25 # DOCTORS: d026_orig, d027_orig, d028_orig, d029_orig
26 ps_num$medici <- compute_factor(ps_num,
27                                c("d026_orig", "d027_orig",
28                                  "d028_orig", "d029_orig"))
29
30 # NURSES: d031_orig, d032_orig, d033_orig, d034_orig
31 ps_num$infermieri <- compute_factor(ps_num,
32                                    c("d031_orig", "d032_orig",
33                                      "d033_orig", "d034_orig"))

```

4.3 Hierarchical Clustering and Cluster Summary

Listing 3: Hierarchical clustering and cluster means

```

1 # Data for clustering: complete cases on the five factor scores
2 cluster_data <- ps_num %>%
3   select(attesa, comfort, assistenza, medici, infermieri) %>%
4   na.omit()
5
6 # Standardize factor scores
7 cluster_scaled <- scale(cluster_data)
8
9 # Euclidean distance matrix
10 dist_mat <- dist(cluster_scaled, method = "euclidean")
11
12 # Ward's hierarchical clustering
13 hc <- hclust(dist_mat, method = "ward.D2")
14
15 # Dendrogram
16 plot(hc,
17       main = "Ward Hierarchical Clustering Dendrogram",
18       xlab = "Observations", sub = "")
19 rect.hclust(hc, k = 3, border = "red")
20
21 # Cut tree into 3 clusters
22 clusters3 <- cutree(hc, k = 3)
23 cluster_data$cluster3 <- clusters3
24
25 # Mean factor scores by cluster
26 cluster_summary <- cluster_data %>%
27   group_by(cluster3) %>%
28   summarise(across(c(attesa, comfort, assistenza, medici, infermieri),
29                     \(x) mean(x, na.rm = TRUE)))
30
31 print(cluster_summary)

```

The console output for the cluster means is:

```

# A tibble: 3 x 6
  cluster3 attesa comfort assistenza medici infermieri
  <int>    <dbl>   <dbl>      <dbl> <dbl>      <dbl>
1       1  0.0327 -0.0527   -0.215 -0.101    -0.122
2       2  0.557   0.432    0.507  0.621     0.658
3       3 -1.21   -0.819   -0.654 -1.06    -1.11

```

4.4 k-means Clustering and Visualisation

Listing 4: k-means clustering and PCA-based cluster plot

```

1 set.seed(123)
2 kmeans_res <- kmeans(cluster_scaled, centers = 3, nstart = 25)
3
4 # Visualisation: fviz_cluster applies PCA for plotting
5 fviz_cluster(kmeans_res,
6               data = cluster_scaled,
7               main = "K-means Clustering on Factor Scores")

```

In the resulting plot, 2, the horizontal axis (Dim1) corresponds to the first principal component, which captures the main gradient of overall satisfaction across the five domains. The vertical axis (Dim2) represents a secondary source of variation, typically less interpretable but useful to visually separate the clusters.

5 Results and Interpretation

5.1 Cluster Profiles

The mean factor scores by cluster (Section 3) allow a straightforward interpretation of the three groups:

- **Cluster 2: Highly Positive Experience**

This group exhibits consistently positive scores across all five dimensions (waiting time, comfort, assistance, doctors, nurses). Patients in this cluster perceive shorter waiting times, better physical conditions, clear communication, and very positive interactions with both medical and nursing staff. This is the “best-experience” group.

- **Cluster 1: Moderately Negative / Average Experience**

This group shows scores around zero or slightly negative, particularly for assistance and staff-related dimensions. Waiting time is close to the overall mean. These patients report a somewhat below-average experience, but not extremely negative. Their dissatisfaction is mainly related to relational and communication aspects rather than delays.

- **Cluster 3: Strongly Negative Experience**

This cluster displays markedly negative scores on all dimensions, with the lowest values for waiting time and staff evaluations. Patients in this group perceive long waiting times, poor comfort, weak communication, and unsatisfactory interactions with doctors and nurses. It represents the most critical profile in terms of ED patient experience.

Both Ward’s hierarchical clustering and k-means identify essentially the same structure: three well-separated groups ordered along a latent “overall satisfaction” dimension.

5.2 Dendrogram and k-means Plot

Figure 1 shows the Ward dendrogram. The height jumps between the last merges support a three-cluster solution, which corresponds to cutting the tree at an appropriate level (red rectangles).

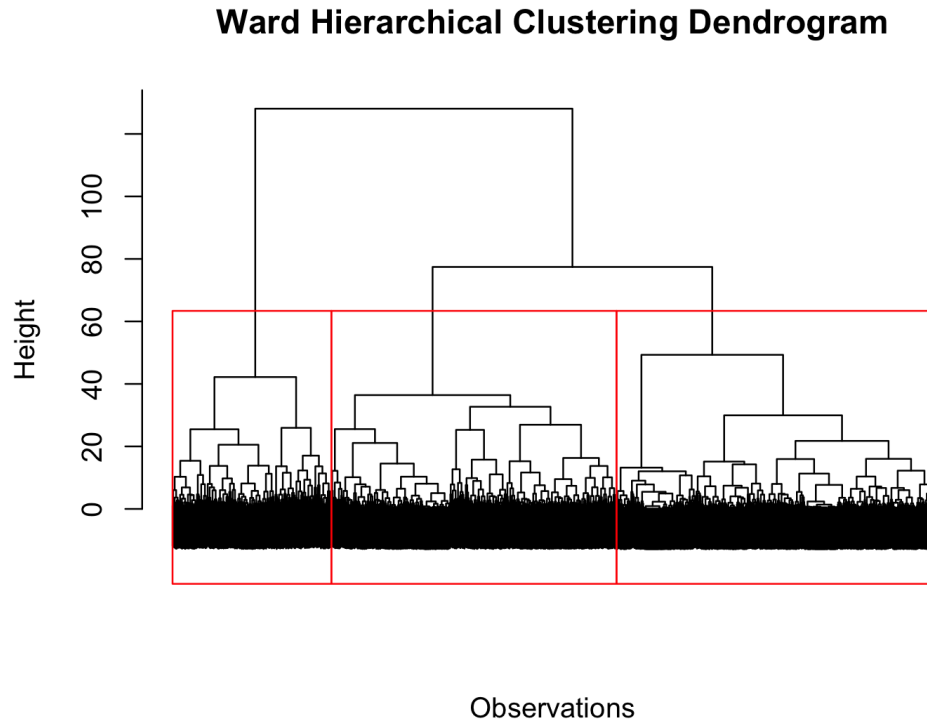


Figure 1: Ward hierarchical clustering dendrogram on the five factor scores.

Figure 2 reports the k-means clustering plot. The axes correspond to the first two principal components of the standardized factor scores. Dim1 typically reflects an overall satisfaction gradient (from negative to positive experiences), while Dim2 captures secondary variation. The three clusters are clearly separated in this low-dimensional representation.

- identifying critical patient segments for targeted interventions;
- illustrating advanced multivariate methods in an applied teaching context.

Students are encouraged to extend the analysis by: adding covariates, exploring additional clustering solutions (e.g. $k = 4$), or linking cluster membership to clinical and organisational outcomes.