# Cluster Correspondence Analysis and Reduced K-Means: A Two-Step Approach to Cluster Low Back Pain Patients

Fengmei Liu, Sucharu Gupta, Cristina Tortora

**Abstract** For the IFCS 2017 data challenge on low back pain (LBP) patients clustering, we used a two-step approach. Two of the challenging characteristics of the data set are the presence of missing values and mixed type variables. After a specific pretreatment, in the first step, we performed domain clustering using cluster correspondence analysis (clusCA). Upon the output variables from each domain, we did the second step, reduced K-means clustering, to get the final clusters of patients. The conclusion section shows the final clustering results and a profile plot of the clusters. Every cluster is highly interpretable and evaluated well with some descriptive variables which are used for measuring the clustering results.

Fengmei Liu
San José State University
✉ lfm.ustb@gmail.com

Sucharu Gupta
San José State University
✉ sucharu7115@gmail.com

Cristina Tortora
San José State University
✉ cristina.tortora@sjsu.edu

# 1 Introduction

The lower back pain data set for the IFCS 2017 data challenge (van Mechelen and Vach (2018)) is characterized by several challenging points; they can be summarized in three main aspects: missing data, high dimensionality, and variables of mixed type. Some of the missing data were not actual missing responses; some of the patients could not answer all the questions. For example, unemployed or retired patients could not answer questions related to the work conditions. Therefore, it was important to find the actual missing values before the imputation phase. Moreover, some of the variables were redundant and should not be used for the analysis; an effective pretreatment was necessary. On the treated data, only a reduced number of clustering techniques can be used. The dimensionality of the data – 121 variables and 928 observations – is too high for most of the classical clustering techniques, like k-means clustering (Hartigan and Wong, 1979) or Gaussian mixture models (Titterington et al, 1985). Also, the data are characterized by mixed type variables; therefore, even the methods for high dimensional data sets, like factor discriminant k-means (Rocci et al, 2011), or mixture of factor analyzers (McLachlan and Peel, 2000), cannot be used. Statistical literature has seen the development of methods for categorical data clustering, see the following for example: Hwang et al (2006); van Buuren and Heiser (1989); D'Enza and Palumbo (2013). However, still few methods work on mixed type data. To face these challenges we propose using a two-step approach. Based on the results from Nielsen et al (2017), we divided the data into domains. In the first step, we performed domain clustering using cluster correspondence analysis (clusCA) (van de Velden et al, 2017). On the continuous variables obtained in step one, we did the second step, reduced K-means clustering (De Soete and Carroll, 1994), to get the final clustering groups for the patients. The obtained clusters are well separated and have meaningful interpretations. The paper is structured as follows:

1. *Data preprocessing* is detailed in Section 2,

2. the *two-step approach* is described in Section 3 and

3. *results and evaluation* are in Section 4.

# 2 Data preprocessing

In the data set there are 121 variables in total. The first 10 of them measure patients' improvement over time; they will not be used for the purpose of clustering but only to describe the results. Table 1 shows the remaining variables categorized per variable type and domain. All the variable indices are the same as shown in the variable file of this data challenge.

**Table 1:** Variable type and domain summary.

| Variable Type | Counts |
| --- | --- |
| Continuous | 8 |
| Dichotomous | 64 |
| Multistate nominal | 9 |
| Ordinal | 30 |
| Trichotomous | 1 |

| Variable Domain | Counts |
| --- | --- |
| Activity | 23 |
| Contextual factors | 16 |
| Pain | 14 |
| Participation | 8 |
| Physical impairment | 24 |
| Psychological | 27 |

## 2.1 Special missing data preprocessing

Most of the variables contain missing values. Before starting the pre-treatment processing, we distinguished missing values by whether people are eligible to respond to the variable questions. The variables from 67 (Fabq60) to 75 (Fabq140) correspond to questions related to work environment (e.g. pain caused by work/accident at work place); therefore, students, unemployed, and pensioner

patients are not eligible to answer the question. We imputed the missing value with the value 0 as a new category. Similarly, for variables 86 to 89, missing values are imputed with 0 as a new category for patients who do not have dominating back pain.

## 2.2 Variable selection

The resulting data set contains only actual missing values. We selected the variables that we included in the analysis according to the following criteria:

- The first 10 variables are not included in the analysis but were used to inform the insights.

- Variables with more than 20% missing data are not included: 91 and 98.

- Summary Scores variables are not included: 85,122.
  Nielsen et al (2016) have pointed out that they got better results by using single items in the questionnaires than using the summary scores, so here we adopted this idea and use single items from the questionnaires.

- Variables that have one categorical level dominating 85% or more (Nielsen et al, 2017) are not included. Variables with a dominating level are: 35, 57, 69, 93, 94, 95, 100, 109, 111, 112, 113, 114, 115 .

The preprocessing resulted in 95 selected variables for the analysis. On top of these 95 variables we did missing value imputation using the random forest method from the *MICE* package (van Buuren and Groothuis-Oudshoorn, 2011) in the R language (R Core Team, 2000).

# 3 Cluster analysis

In this section we describe the techniques used in our two-step approach: cluster correspondence analysis and reduced K-means. Cluster correspondence analysis (clusCA) is a method for joining dimension reduction and cluster analysis for

categorical data proposed by van de Velden et al (2017). It is based on multiple correspondence analysis (MCA), a well-known technique for categorical variable reduction (Greenacre and Pardo, 2006); it extends correspondence analysis (CA) (Greenacre, 1984) to the case of multivariate data sets. MCA and CA, however, only perform dimension reduction. In the recent literature, several methods that jointly or sequentially perform dimensional reduction and cluster analysis have been proposed; among them are, GROUPALS (van Buuren and Heiser, 1989), MCA K-means (Hwang et al, 2006), and i-FCB (D'Enza and Palumbo, 2013). van de Velden et al (2017) show that clustCA performs better than the cited methods. Let us define with $K$ the number of clusters, $n$ the number of observations, $d$ the dimension of the reduced subspace, and $q$ the number of categorical variables. Each categorical variable has $p_j$ modalities, with $j = 1 \ldots q$, and $Q = \sum_{j=1}^{q} p_j$. The following formula shows the clusCA optimization criterion:

$$min\ \phi_{clusca}\ '(\mathbf{Z}_K, \mathbf{G}) = \|\sqrt{\frac{n}{q}}\mathbf{M}\mathbf{Z}\mathbf{D}_z^{-\frac{1}{2}}\mathbf{B}^* - \mathbf{Z}_K\mathbf{G}\|^2 \qquad (1)$$

where $\mathbf{B}^* = \frac{1}{\sqrt{nq}}\mathbf{D}_z^{1/2}\mathbf{B}$, $\mathbf{D}_z = diag(\mathbf{Z}'\mathbf{Z})$ with

| | |
|---|---|
| $\mathbf{M}$ | centering operator of $\mathbf{Z}$, |
| $\mathbf{Z}$: $n \times Q$ | binary matrix, where $\mathbf{Z}_j$ is an $n \times p_j$ indicator matrix for the j-th categorical variable with $j = 1 \ldots q$, |
| $\mathbf{B} = [\mathbf{B}_1', \mathbf{B}_2', ..., \mathbf{B}_q']'$ | as the $Q \times d$ matrix of orthonormal loadings also interpreted as category quantifications, |
| $\mathbf{Z}_K$: $n \times K$ | binary matrix indicating cluster memberships, and |
| $\mathbf{G}$: $K \times d$ | cluster centroid matrix. |

The package also contains good visualization tools for all the methods implemented in the package. Reduced K-means is a popular subspace clustering method; it is designed to maximize the between-cluster deviance. Let us define with $\mathbf{X}$ an $n \times p$ continuous data matrix, with $n$ number of observations and $p$ number of variables. Reduced K-means alternatively looks for a reduced subspace of dimension $q < p$ and a data partition in $K$ clusters. Specifically, let's define with $\mathbf{A}$ a $p \times q$ component weights matrix for variables, with $\mathbf{U}$ an

$n \times K$ membership matrix defining a partition of units into K clusters, where $u_{ik} = 1$ if the $i^{th}$ object belongs to cluster k, $u_{ik} = 0$ otherwise, and **M** a $K \times q$ centroid matrix in the reduced space, where $m_{kj}$ is the centroid value of the $j^{th}$ component obtained on the $k^{th}$ cluster. Reduced K-means (De Soete and Carroll, 1994) looks for **A**, **M**, **U** that minimize

$$F_{rkm}(\mathbf{A}, \mathbf{M}, \mathbf{U}) = ||\mathbf{X} - \mathbf{UMA}'||^2 . \tag{2}$$

Reduced K-means returns a partition of the data in $K$ clusters and the projection of the data in the reduced subspace. The code can be found in the R package *clustrd* (Markos et al, 2013).

## 3.1 Step 1 - cluster correspondence analysis

The strategy is to treat all the variables as categorical. For the continuous variables, we transform them to categorical variables based on quantiles. These variables are

1. Age,
2. Bhoej0 (Height),
3. Vasl0 (LBP intensity),
4. Okon0 (Able to decrease pain),
5. Obeh0 (Treatment not essential),
6. Htil0 (Self-rated general health) and
7. bmi.

After transforming all the continuous variables to categorical, we used clusCA based on the six domains that are suggested by Nielsen et al (2017). The goal is to extract the main components that represent each domain. This step was done using a tuning process, specifying the groups to be 3:12, and reduced dimensions to be 2:9, and using the criterion of average silhouette width to choose the best reduced dimensions and number of clusters.

To interpret and use the results, we utilize the variable component information rather than the clustering information. First of all, we looked into the main biplot, and summarized the main variables that contribute to the positive (H)

and negative (L) direction of selected components. Second, the feature plots of each cluster further showed and helped the variable interpretation from the main biplot. They show the top 20 variables in each cluster, and positive sign means the variable frequency is above average in the cluster, while negative sign means the variable frequency is below average in the cluster.

Now we present the output and analysis for the six domains. As an illustration of the notation, "HX1" means the positive direction of the first component from the contextual factor domain. "LX4" is the negative direction of the second component from the activity domain.

**Contextual factor domain**
The output has four clusters and two components. Details of each component are shown in appendix figure 4:

HX1      male, age 34 to 43, tall, no chronic disease, full-time work, low education,

LX1      female, old, short, has musculoskel/other chronic disease, retired, high education,

HX2      female, young, tall, no other chronic disease, full time work, high education and

LX2      male, old, average height, has musculoskel/other chronic disease.

**Activity domain**
The output has three clusters and two components. Details of each component are shown in appendix figure 5:

HX3      home activity slowly, self dress slowly,

LX3      home activity normal, self dress normal,

HX4      walk short distance, stand for short time, self-dress normal and

LX4      walk normal, stand normal, self-dress slowly.

**Pain domain**

The output has three clusters and two components. Details of each component are shown in appendix figure 6:

HX5     pain has spread down to legs, leg pain is intense,

LX5     pain has not spread down to legs, no leg pain,

HX6     LBP pain intensity is high, this pain episode last short and

LX6     LBP pain intensity is low, this pain episode last long.

**Participation domain**

The output has four clusters and two components. Details of each component are shown in appendix figure 7:

HX7     work does not make pain worse, physical work load sitting/walking,

LX7     work does make pain worse, heavy physical work,

HX8     work is heavy, more sick leave time and

LX8     unsure if work is heavy, less sick leave time.

**Physical impairment domain**

The output has three clusters and two components. Details of each component are shown in appendix figure 8:

HX9     no pain on AROM,

LX9     leg pain on AROM test,

HX10    negative on SI-joint tests, no pain on palpation and

LX10    positive on SI-joint tests, back pain on AROM test.

**Psychological impairment domain**

The output has three clusters and two components. Details of each component are shown in appendix figure 9:

HX11    good mood and feel energetic, good conscience, good sleep,

LX11    bad mood and feel less energetic, bad conscience, bad sleep,

HX12    not loose interest in daily activities, psychologically believe should do activities and

LX12    lose interest in daily activities, psychologically believe should not do activities.

## 3.2 Step 2 - Reduced K-means

The new component output of clusCA is now in the Euclidean space. Therefore, the inputs for reduced K-means are the combined output of X1 to X12 from the above clusCA results of six domains. Here, we used all the 12 dimensions instead of any dimension reduction. To find the number of clusters the commonly used indices are: within variance, silhouette width Rousseeuw (1987), and the Calinski-Harabasz (CH) index (Desgraupes, 2013). Each index measures the quality of a partition using different criteria, specifically a small within variance, or equivalently a large between variance, guarantees that the elements in a cluster are closer to the elements in the same cluster than to the elements in the other clusters. Similarly, the CH index is based on a ratio of between and within clusters variance. A high CH indicates a good partition. The silhouette value measures how similar an object is to its own cluster compared to other clusters using a distance measure. Therefore, the optimal number of clusters is the one associated with a low within variance, and high silhouette width and CH index. We also used the self tuning RKM function in R package *clustrd* to choose the number of clusters, and, given the dimensions, we specified the groups to be 2:20, and dimensions to be 12. From the plots of the three measures in figure 1, we can tell that a number of groups in the range 5 to 8 optimizes the mentioned criteria. We created multiple partitions of the data fixing the number of groups in the selected range and finalized with eight groups, which has the largest silhouette value, high CH index, around the within variance elbow point, and also good interpretations of each cluster.

**Figure 1:** Within variance, silhouette width, and CH index for different numbers of clusters. Note: to make the graph compact, we put all three measures in one plot, with rescaling within variance/10000, Silhouette width*3, and CH index/300.

# 4 Results

The clustering results are shown in the profile plot in figure 2a and the interpretations of all clusters are given in Table 2 on pages 12 to 13.



(a) Final clustering profile plot.



(b) 3D visualization using Rtsne.

**Figure 2:** The resulting graphs..

Using the dimension reduction method Rtsne (van der Maaten, 2014), we can visualize the clusters in a lower dimensional space. Figure 2b shows a 3D visualization of the clusters; they are well separated.



(a) LBP intensity.  (b) rmprop.  (c) gen.

**Figure 3:** Clustering evaluation with first 10 variables..

By using the first 10 variables, we compare the clustering results with LBP Intensity (vasl2w, vasl3m, vasl12m), figure 3a, RM summary score (rmprop2w, rmprop3m, rmprop12m), figure 3b, and global perceived improvements (gen2w, gen3m, gen12m), figure 3c.

There are two important evaluation qualities we can highlight from here: Clusters are interpretable and meaningfully match with the evaluation variables, and clusters are well separated. Specifically, focusing on the interpretability, in the left plot of LBP intensity in figure 3, clusters 1 and 2 are the lowest; in our results, cluster 1 is composed of middle-aged men who have intense LBP with a short duration, so the pain healed quickly and dropped down fast. Cluster 2 is a group of young women with light LBP; therefore, the intensity drop is not large but it belongs to the lowest generally. In contrast, cluster 6 is a group of older people with pain and positive SI-joint and AROM tests, and cluster 8 is a group that has pain which has spread to legs intensely. So their LBP pain is on average higher than other groups at all the time points. Similarly, for the middle plot of RM summary score, clusters 1 and 2 are shown as the lowest and clusters 8 and 6 are higher than other groups.

Overall, the interpretation of the clustering results matches with the evaluation. Moreover, all the clusters are well separated at two weeks, three months, and 12 months at each variable. This quality indicates a good separation of the observations.

**Table 2:** Clustering result and interpretation (1/2).

| Cluster | Percentage* | Main Features | Interpretation (features are separated by ";") |
|---|---|---|---|
| C1 | 23.6% | HX1, HX4, LX5, HX6, HX7, HX9, HX10, HX11 | Mostly male, age 34 to 43, no chronic disease, full-time work; Walk short distance, stand for short time, self-dress normal; Pain has not spread down to legs; LBP pain intensity is high; This pain episode is short; Work does not make pain worse, physical workload sitting/walking; no pain on AROM test; Negative on SI-joint tests, no pain on palpation; Good mood and feel energetic, good sleep. |
| C2 | 17.9% | HX2, LX3, LX4, LX5, LX6, HX7, HX11, HX12 | Mostly female, young, no other chronic disease, full-time work, high education; Normal home activity; Walks normally, stands normally; Pain has not spread down to legs, no leg pain; LBP pain intensity is low, this pain episode is long; Work does not make pain worse, physical workload is sitting/walking; Good mood and feel energetic, good conscience, good sleep; Not loose interest in daily activities, psychologically believe should do activity. |
| C3 | 14.9% | HX2, HX4, HX5, LX6, HX10 | Walk short distance, stand for short time, self-dress normal; Pain has spread down to legs, leg pain is intense; LBP pain intensity is low, this pain episode is long; Negative on SI-joint tests, no pain on palpation. |
| C4 | 12.7% | LX1, HX3, HX6, LX10 | Mostly female, old, short, has musculoskel/other chronic disease, retired, high education; Home activity slowly, self dress slowly; LBP pain intensity is high, this pain episode is short; Positive on SI-joint tests, back pain on AROM test. |

* percentage means the size of each cluster, i.e. number of patients in each cluster divided by all the studied population

**Table 2:** Clustering result and interpretation (2/2).

| Cluster | Percentage* | Main Features | Interpretation (features are separated by ";") |
|---------|-------------|---------------|------------------------------------------------|
| C5 | 12.3% | HX1, LX4, HX6, LX7, HX8, LX12 | Mostly male, age 34 to 43, tall, no chronic disease, full-time work, low education; Walks normally, stands normally, self-dress slowly; LBP pain intensity is high, this pain episode last short; Work makes pain worse, heavy physical work; Work is heavy, more sick leave time; Loose interest in daily activities, psychologically believe should not do activities. |
| C6 | 8.1% | LX1, LX2, LX10, LX11, HX12 | Old people, has musculoskel/other chronic disease, retired; Positive on SI-joint test, back pain on AROM test; Bad mood and feel less energetic, bad conscience, bad sleep; Not loose interest in daily activities, psychologically believe should do activity. |
| C7 | 6.7% | LX2, LX7, LX8 | Mostly male, old (52 to 66), average height, has musculoskel/other chronic disease, (low education); Work makes pain worse, heavy physical work; Unsure if work is heavy, less sick leave time. |
| C8 | 3.9% | HX5, LX9, HX10 | Pain has spread down to legs, leg pain is intense; Leg pain on AROM test; Negative on SI-joint tests, no pain on palpation. |

* percentage means the size of each cluster, i.e. number of patients in each cluster divided by all the studied population

# Appendix



**Figure 4:** clusCA output of contextual factor domain.

**Figure 4:** clusCA output of contextual factor domain[1].

In the contextual factor domain, as shown in the upper left plot, the clustering is best fitted as 3 clusters and 2 components, using the criterion of average silhouette width. The upper right plot shows that the first cluster has 38.1% of all the patients, and each bar represents the contribution of the top 20 variables and categories. Positive sign means the variable frequency is above average in the cluster, while negative sign means the variable frequency is below average in the cluster. Similarly, the bottom three plots from left to right show the cluster sizes(shown as percentages) and variable contributions in the other three clusters.



**Figure 5:** clusCA output of activity domain[2].

**Figure 5:** clusCA output of activity domain[3].

**Figure 5:** clusCA output of activity domain[4].

The interpretation of figure 5 is similar to figure 1. The cluster plots (the upper right and bottom two plots) only show the top 20 variables' contribution. The grey background subplots shown in them are an overview of all variables' contribution. Similar interpretation applies to figures 6 - 9.

**Figure 6:** clusCA output of pain domain.

**Figure 6:** clusCA output of pain domain.

**Figure 7:** clusCA output of participation domain.

**Figure 7:** clusCA output of participation domain.

**Figure 7:** clusCA output of participation domain.



**Figure 8:** clusCA output of physical impairment domain.

**Figure 8:** clusCA output of physical impairment domain.

**Figure 8:** clusCA output of physical impairment domain.



**Figure 9:** clusCA output of psychological impairment domain.

**Figure 9:** clusCA output of psychological impairment domain.

**Figure 9:** clusCA output of psychological impairment domain.

# References

van Buuren S, Groothuis-Oudshoorn K (2011) Mice: Multivariate Imputation by Chained Equations in R. Journal of Statistical Software 45(3):1–67, DOI 10.18637/jss.v045.i03

van Buuren S, Heiser WJ (1989) Clustering n objects into k groups under optimal scaling of variables. Psychometrika 54(4):699–706, DOI 10.1007/BF02296404

De Soete G, Carroll JD (1994) K-means clustering in a low-dimensional Euclidean space. In: New approaches in classification and data analysis. Studies in Classification, Data Analysis, and Knowledge Organization, Diday E, Lechevallier Y, Schader M, Bertrand P, Burtschy B (eds), 1st edn., Springer, Berlin, Heidelberg (Germany), p. 212–219, DOI 10.1007/978-3-642-51175-2_24

D'Enza AI, Palumbo F (2013) Iterative factor clustering of binary data. Computational Statistics 28(2):789–807, DOI 10.1007/s00180-012-0329-x

Desgraupes B (2013) Clustering indices. Journal of University of Paris Ouest-Lab Modal'X 1:1–34

Greenacre M (1984) Theory and application of correspondence analysis. Academic Press, London (UK). ISBN: 978-0-122990-50-2

Greenacre M, Pardo R (2006) Multiple correspondence analysis of subsets of response categories. In: Multiple Correspondence Analysis and Related Methods, Greenacre M, Blasius J (eds), Chapman & Hall/CRC, Boca Raton (USA), p. 197–217, DOI 10.2139/ssrn.847647

Hartigan JA, Wong MA (1979) Algorithm AS 136: A k-means clustering algorithm. Journal of the Royal Statistical Society. Series C (Applied Statistics) 28(1):100–108, JSTOR / Wiley for the Royal Statistical Society, DOI 10.2307/2346830

Hwang H, Montréal H, Dillon WR, Takane Y (2006) An Extension of Multiple Correspondence Analysis for Identifying Heterogeneous Subgroups of Respondents. Psychometrika 71(1):161–171, DOI 10.1007/s11336-004-1173-x

van der Maaten L (2014) Accelerating t-SNE using tree-based algorithms. Journal of Machine Learning Research 15(1):3221–3245, URL http://www.jmlr.org/papers/volume15/vandermaaten14a/vandermaaten14a.pdf

Markos A, D'Enza AI, van de Velden M (2013) Clustrd: Methods for joint dimension reduction and clustering. R package version 0.1.2, URL https://rdrr.io/cran/clustrd/

McLachlan GJ, Peel D (2000) Mixtures of factor analyzers. In: Proceedings of the Seventh International Conference on Machine Learning, Morgan Kaufmann, San Francisco (USA), p. 599–606

van Mechelen I, Vach W (2018) Cluster analyses of a target data set in the IFCS cluster benchmark data repository: Introduction to the special issue. Archives of Data Science, Series B 1(1):1–12, DOI 10.5445/KSP/1000085952/01

Nielsen AM, Vach W, Kent P, Hestbaek L, Kongsted A (2016) Using existing questionnaires in latent class analysis: should we use summary scores or single items

as input? A methodological study using a cohort of patients with low back pain. In: Clinical epidemiology, Sørensen HT, Cohen E, Ehrenstein V, Pedersen L, Petersen I, Benchimol E, Briko N, Brookhart MA, Browner WS, Duhaut P, Field JK, Jick S, Jin L, Kieler H, Ness RB, Shah NH, Smeeth L, Trifiro G, Vlassov VV, Wang J, Weiss NS, Zhan S (eds), vol. 8, Dove Press, p. 73–89, DOI 10.2147/CLEP

Nielsen AM, Kent P, Hestbaek L, Vach W, Kongsted A (2017) Identifying subgroups of patients using latent class analysis: Should we use a single-stage or a two-stage approach? A methodological study using a cohort of patients with low back pain. BMC Musculoskeletal Disorders 18(1):57, BioMed Central, DOI 10.1186/s12891-017-1411-x

R Core Team (2000) R language definition. R Foundation for Statistical Computing, URL https://cran.r-project.org/doc/manuals/r-release/R-lang.html

Rocci R, Gattone SA, Vichi M (2011) A New Dimension Reduction Method: Factor Discriminant K-means. Journal of Classification 28(2):210–226, DOI 10.1007/s00357-011-9085-9

Rousseeuw PJ (1987) Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics 20:53–65, DOI 10.1016/0377-0427(87)90125-7

Titterington DM, Smith AFM, Makov UE (1985) Statistical analysis of finite mixture distributions. Wiley series in probability and mathematical statistics, Wiley, Chichester (UK), New York (USA). ISBN: 04-7190-763-4

van de Velden M, D'Enza AI, Palumbo F (2017) Cluster Correspondence Analysis. Psychometrika 82(1):158–185, Springer US, DOI 10.1007/s11336-016-9514-0