# Assignment 2

Michael Darmanis (7115152200004)          Vasilios Venieris (7115152200017)
mdarm@di.uoa.gr                          vvenieris@di.uoa.gr

January 24, 2023

## 1   Introduction

The main goal of this report[1] is to provide recommendations to the CEO of HELP International on how to allocate the organisation's resources in a strategic and effective manner, based on the identification of countries in need of development aid. This will be achieved through the implementation of a clustering analysis on a country data set using MATLAB®. The analysis will consist of several stages, including the exploration of the data's characteristics and patterns, the selection and transformation of relevant features, the selection of appropriate clustering algorithms, the execution of these algorithms, and the characterization of the resulting clusters.

The initial stage of the analysis involves the examination of the individual features of the data, including their data type and range of values, as well as the distribution of these values. The relationships between different features will also be considered. Following this, relevant features will be selected and transformed in order to make their values comparable. Based on the characteristics of the data and the desired properties of the clusters, the appropriate clustering algorithms will then be chosen. These algorithms will then be executed with different parameter-values in order to identify persistent clusters. Finally, the clusters will be characterised based on the values of the features within each cluster.
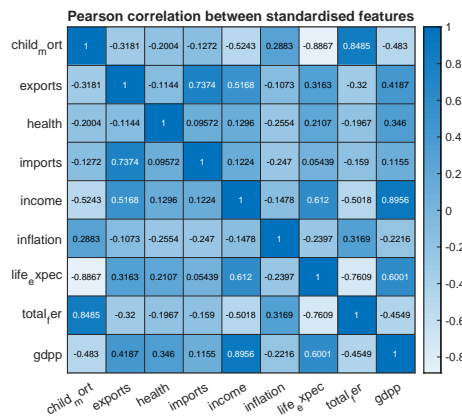
## 2   "Feeling the data"

In this part of the analysis, characteristics of each individual feature in the dataset were examined in order to gain a better understanding of its nature. This involved determining the data type and range of values for each feature, as well as generating histograms to visualize the distribution of values. The mean and standard deviation were also calculated for each feature.

```
Feature      Type      Range                            Mean         Std Dev
----------   --------  -----------------------------    ----------   -----------
 child_mort  float     [      2.6000,     208.0000]        38.2701       40.3289
 exports     float     [      0.1090,     200.0000]        41.1090       27.4120
 health      float     [      1.8100,      17.9000]         6.8157        2.7468
 imports     float     [      0.0659,     174.0000]        46.8902       24.2096
 income      float     [    609.0000,  125000.0000]     17144.6886    19278.0677
 inflation   float     [     -4.2100,     104.0000]         7.7818       10.5707
 life_expec  float     [     32.1000,      82.8000]        70.5557        8.8932
 total_fer   float     [      1.1500,       7.4900]         2.9480        1.5138
 gdpp        float     [    231.0000,  105000.0000]     12964.1557    18328.7048
```
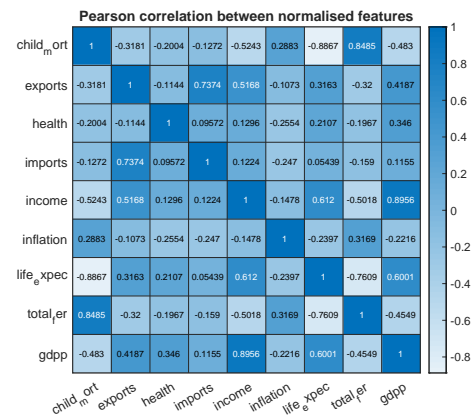
The linear dependence between each feature and all the others was calculated using the Pearson correlation coefficient in order to gain insight into the relationships between different features in the data set. Additionally, standard score normalization and min-max feature scaling normalization were performed. The linear dependence between the transformed features were calculated in each case.

It is evident in Figure 1 that linear relationships between features are not affected by the transfor-
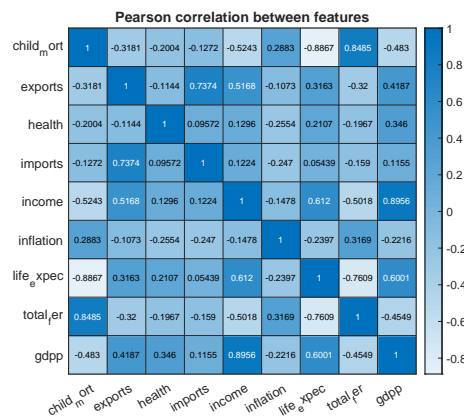
---

[1]as stated at https://www.kaggle.com/rohan0301/unsupervised-learning-on-country-data

---

(a) Standardised features.



(b) Min-max normalised features.



(c) Raw features.

Figure 1: Correlation matrices between features.

mations. Pearson's correlation measures the linear component of association so it comes to no surprise that linear transformations of data (like mim-max normalisation and standardisation) did not affect the correlations between the features.

In Figure 2 it can readily be observed that the feature `life_expec` exhibits a negative skew in its distribution, while the `health` exhibits a normal distribution. On the other hand, all other measured quantities show a positive skew in their distributions. It should be noted that the distribution of the categorical feature "country", which consists solely of text data and exhibits the same number of unique values as the total length of the data set, was not analysed since it makes little sense to do so.
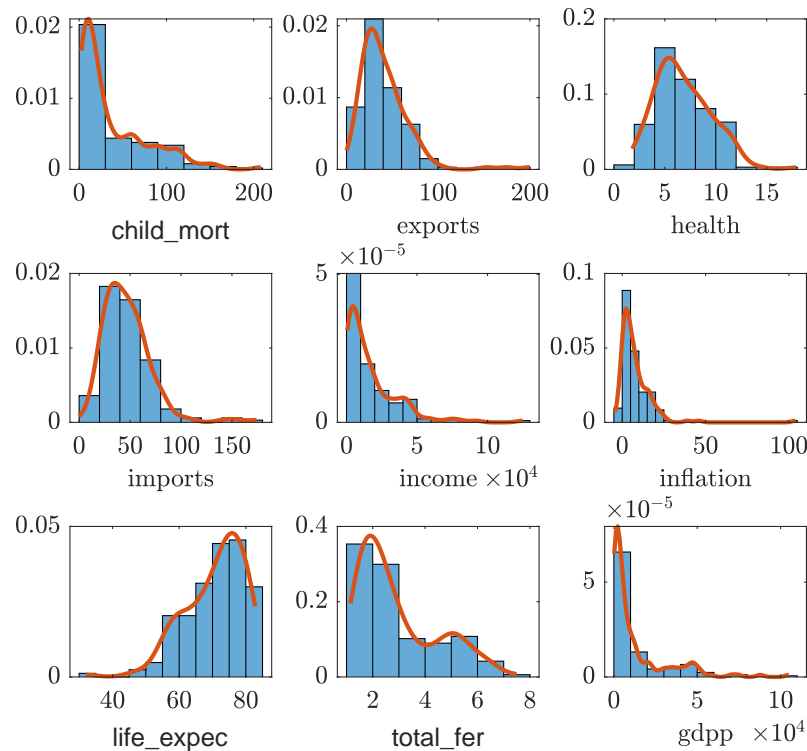
Figure 2: Histograms and distributions of features

The analysis of the dataset also reveals several relationships between features. The relationship between child mortality and economic conditions is of particular significance, as the data indicates that child mortality tends to increase as income, gross domestic product (GDP), and exports decrease. Inflation also appears to have a negative impact on child mortality. This suggests that economic factors, such as income and GDP, may be influential in determining child mortality rates. The relationship between exports and other economic indicators is also noteworthy, as an increase in exports tends to lead to an increase in GDP, income, and imports, implying that exports may be a key contributor to economic growth.

In addition, the data suggests that spending on health has a positive effect on life expectancy and a negative effect on child mortality. Higher levels of income and GDP are correlated with higher life expectancy and lower child mortality, indicating a possible relationship between these factors and spending on health. Furthermore, high levels of inflation appear to be detrimental to economic conditions, as they have a negative effect on various economic indicators, including income, GDP, and total fertility rate. Finally, the data suggests that higher life expectancy is correlated with lower total fertility rates, a relationship that may be influenced by factors such as GDP and spending on health.

# 3    Feature selection/transformation

Based on the aformentioned data relationships, it is clear that some features are closely related to specific categories, namely: health, trade, and finance. Therefore, features will be grouped up into these categories and then normalised using a min-max scheme. The three categories of features in the dataset are: health (child mortality, health, life expectancy, total fertility rate), trade (imports, exports), and finance (income, inflation, GDP).
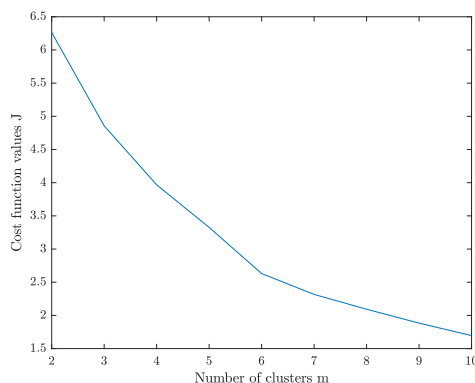
By combining multiple features into fewer ones facilitates the comprehension of relationships between different features and provides a more intuitive understanding of the data. It also captures broader or more general relationships in the data, and it may be possible to improve the generalisability of the clustering results to new, unseen data.

Additionally, by creating new features that capture more relevant or subtle relationships in the data, it may be possible to improve the performance of the clustering algorithm itself.
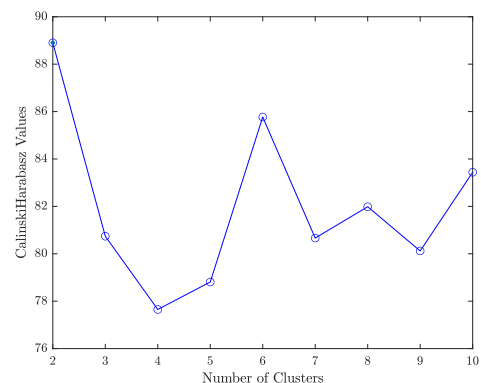
# 4    Selection and execution of the clustering algorithms

After normalising and combining & renormalising the data, the k-means clustering algorithm was used for performing the clustering analysis in order to effectively group countries in terms of their needs. The k-means algorithm is a popular choice for clustering due to its simplicity and efficiency, making it well-suited for large datasets (not that ours is, at this point anyway). Additionally, the algorithm produces clear and distinct clusters, which can be beneficial when attempting to partition the countries into groups. Overall, the use of k-means clustering in this context can provide valuable insights and aid in the strategic allocation of resources.

Determining the appropriate number of initial clusters in a k-means clustering operation can be achieved through the use of various measures such as the "elbow method" or the variance-ratio criterion[1]. The "elbow method" for approximating the correct value of k is to run the algorithm for increasing values of k, until there is a minimal decrease in the chosen measure of cluster cohesion (the cost function in our case) between two values[5]. Alternative measures such as the Bayesian Information Criterion and Gap statistics are also suggested as preferable options[7]. It is important to note that the "elbow method", which involves identifying a point of inflection in a plot of the measure of cluster cohesion versus k, has been criticized in the literature [6, 4] and should, generally speaking, be avoided.



(a) "Elbow" method.                    (b) Variance-ratio criterion.

Figure 3: Picking the right number of initial clusters k.

Both the "elbow method" (see Figure 3a) and the variance-ratio criterion (see Figure 3b) agreed that the hidden structure underlying the data was expressed in three clusters. The the k-means algorithm was initialised using three clusters.

In order to pick initial points that have a good chance of being in different clusters, k-means centroids were randomly initialized from Gaussian noise of the data as described in Coates et al. [2]. The algorithm

was then applied 100 times to the dataset, and the results with the smallest cost function were selected. The resulting clusters are shown in Figure 4.
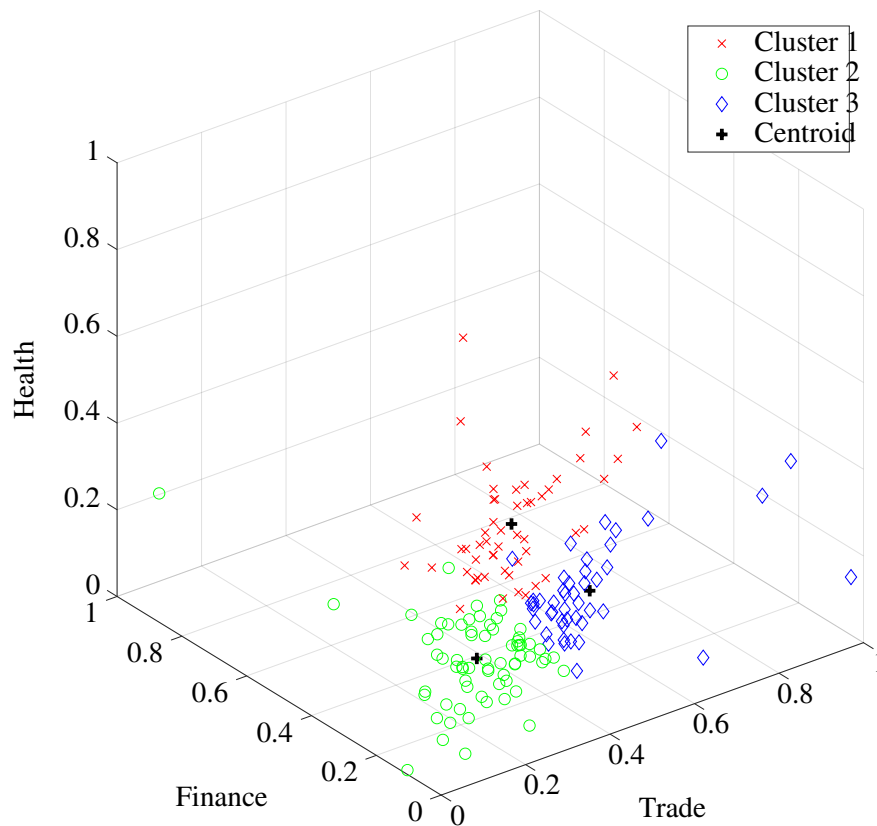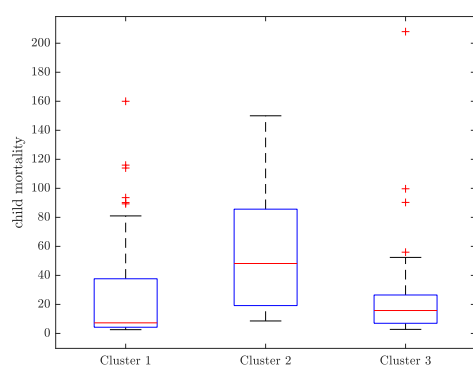


Figure 4: Clusters

It is important to note that running the algorithm multiple times helps to make the k-means more tolerant to outliers and results in the final centroid leaning towards the denser are of points. If that were not the case, then the algorithm would have been sensitive to outliers, and a variant such as k-medians ought have been used instead as described in Theodoridis and Koutroumbas [8].

## 5   Characterisation of clusters

As seen in Section 2, a strong correlation exists between low income and high child mortality. This relationship is widely acknowledged within the fields of economics and public health (see Appendix A). Low income is frequently considered an indicator of economic underdevelopment, and high child mortality reflects the overall health and well-being of a population.

Based on the strong correlation between low income and high child mortality, it is safe to assume that countries with low income and high child mortality rates are likely to have less developed economies and weaker healthcare systems, thereby being in need of funding. Figure 5 showcases the income and child mortality rates with respect to labelled clusters. By examining Figures 5a and 5b, we can make an initial assumption and say that the countries belonging to Cluster 1 are the ones in more need of development aid.

(a) Child mortality.



(b) Net income.

Figure 5: Death of children under 5 years of age per 1000 live births and net income per person of all the countries within each Cluster.

By examining the countries within Cluster 1, we can observe that even developed countries such as the United Kingdom are included in this cluster. This suggests that the clustering process and feature engineering[2] should be re-evaluated in order to accurately identify which countries are in need of funding.

```
        Cluster 1               Cluster 2               Cluster 3
       ----------              ----------              ----------

      Afghanistan                 Angola                   Benin
     Burkina Faso                Burundi                 Albania
          Algeria    Antigua and Barbuda               Argentina
          Armenia              Australia                 Austria
          Bahrain                Belgium                  Brunei
                .                      .                       .
                .                      .                       .
                .                      .                       .
         Tanzania            Timor-Leste                    Togo
           Uganda                 Zambia                 Uruguay
       Uzbekistan                Vanuatu                 Vietnam
            Yemen    Switzerland United Arab Emirates
   United Kingdom          United States               Venezuela
```

Prima facie evidence suggests that there could be a case for further investigation, to establish whether or not further clustering analysis should be put in hand. Nevertheless, it should be stressed that, an unsupervised analysis is limited and relevant facts could be difficult to establish with any degree of certainty.

# References

[1]  Tadeusz Calinski. "A dendrite method for cluster analysis". In: *Communication in statistics* 3 (1974), pp. 1–27.

[2]  Adam Coates and Andrew Y Ng. "Learning feature representations with k-means". In: *Neural networks: Tricks of the trade*. Springer, 2012, pp. 561–580. URL: https://link.springer.com/chapter/10.1007/978-3-642-35289-8_30.

---

[2]an alternative is to discard highly correlated data altogether and keep only low-valued correlations; such a case is seen in [3]

[3]  Tanmay Deshpande. *Clustering: PCA | K-Means - DBSCAN - Hierarchical |*. 2022. URL: https://www.kaggle.com/code/tanmay111999/clustering-pca-k-means-dbscan-hierarchical/notebook.

[4]  David J Ketchen and Christopher L Shook. "The application of cluster analysis in strategic management research: an analysis and critique". In: *Strategic management journal* 17.6 (1996), pp. 441–458.

[5]  Jure Leskovec, Anand Rajaraman, and Jeffrey David Ullman. "Clustering". In: *Mining of massive data sets*. Cambridge university press, 2020. Chap. 7, pp. 240–280. URL: http://mmds.org/#ver30.

[6]  Glenn W Milligan and Martha C Cooper. "An examination of procedures for determining the number of clusters in a data set". In: *Psychometrika* 50.2 (1985), pp. 159–179.

[7]  Erich Schubert. "Stop using the elbow criterion for k-means and how to choose the number of clusters instead". In: *arXiv preprint arXiv:2212.12189* (2022).

[8]  Sergios Theodoridis et al. "Clustering". In: *Introduction to Pattern Recognition: A Matlab Approach*. Academic Press, 2010. Chap. 7, pp. 159–208.

## Appendix A

According to the World Bank, low income is defined as a gross national income (GNI) per capita of less than $1,035 per year, while high child mortality is defined as a mortality rate of more than 43 deaths per 1,000 live births. Countries with low income and high child mortality rates tend to have less developed economies and weaker healthcare systems, which can contribute to higher mortality rates among children.

Many factors can contribute to low income and high child mortality in a country, including poverty, lack of access to education and healthcare, lack of infrastructure and resources, and political instability. Developing countries often face these challenges to a greater extent than developed countries, which can lead to higher rates of child mortality and lower levels of economic development.

Several sources discuss the relationship between low income and high child mortality, including:

- The World Health Organization (WHO) states that "poverty is the single greatest threat to child survival" and that "most of the 10.6 million deaths among children under five occur in developing countries." (https://www.who.int/child_adolescent_health/topics/poverty/en/)

- The World Bank's World Development Indicators database provides data on income levels and child mortality rates for countries around the world. (https://data.worldbank.org/indicator/SP.DYN.IMRT.IN)

- The United Nations Children's Fund (UNICEF) reports that "more than half of all child deaths occur in just five countries – India, Nigeria, Pakistan, Democratic Republic of the Congo and China – with India alone accounting for about one third." (https://www.unicef.org/sowc2014/numbers/)

## Appendix B

Listing 1: Clustering analysis script

```
1  %% Clustering Algorithms, Homework 2
2  %  Clustering analysis on country data in order to determine which
3  %  group of countries are in need of financial aid.
4  %
5  %  This script makes use of the following provided functions:
6  %
```

```matlab
7  %        rand_data_init.m
8  %        k_means.m
9  %
10 %   Two further functions, written by the authors, were also used
11 %   for editing plot variables, namely:
12 %
13 %        PlotDimensions.m
14 %        ChangeInterpreter.m
15 %

17 %% Initialisation
18 clear; close all; clc

20 %% ================= Part 1: Feeling the data
       =========================

22 % Import the CSV file
23 data = readtable('Country-data.csv');

25 % Extract the labels from the first column
26 countryNames = data{:, 1};

28 % Extract the column names and keep only feature labels
29 featureNames = data.Properties.VariableNames;
30 featureNames = featureNames(1, 2:end);

32 % Convert the table-data to an array
33 data = table2array(data(:, 2:end));

35 % Determine the dimensions of the data set
36 [numRows, numCols] = size(data);

38 % Preallocate cell array for feature type
39 featureType = cell(1, 9);

41 % Print a header row for the table
42 fprintf('%-13s %-9.5s %-30s %-11s %-11s\n', 'Feature',...
43     'Type', 'Range', 'Mean', 'Std Dev');
44 fprintf('%-13s %-8s  %-30s %-11s %-11s\n',...
45     '-----------', '--------',...
46     '------------------------------', '-----------', '-----------');

48 % Create a grid of subplots, with one subplot for each feature
49 figure(1);
50 subplot(3, 3, 1:numCols);

52 for i = 1:numCols
53     uniqueVal = unique(data(:, i));
54     if isstring(uniqueVal)
55         featureType(1, i) = {'categorical'};
56     elseif isinteger(uniqueVal)
57         featureType(1, i) = {'integer'};
```

```matlab
58          elseif isfloat(uniqueVal)
59              featureType(1, i) = {'float'};
60          end
61
62          % Determine the range of values, mean, and standard deviation
63          % for the current column
64          minVal  = min(data(:, i));
65          maxVal  = max(data(:, i));
66          meanVal = mean(data(:, i));
67          stdVal  = std(data(:, i));
68
69          % Print a row for the current column
70          fprintf(' %-11s  %-8s  [%12.4f, %12.4f] %12.4f %12.4f\n',...
71              featureNames{i}, featureType{i}, minVal, maxVal,...
72              meanVal, stdVal);
73
74          % Select the subplot for the current feature
75          subplot(3, 3, i);
76
77          % Extract the data for the current feature
78          featureData = data(:, i);
79
80          % Create a histogram for the current feature
81          histogram(featureData, 'Normalization', 'pdf');
82
83          % Allow the distribution plot to be superimposed on the
84              histogram
85          hold on;
85          x = linspace(min(featureData), max(featureData), 100);
86          y = ksdensity(featureData, x);
87          plot(x, y, 'LineWidth', 2);
88          xlabel(featureNames(i), 'Interpreter', 'none');
89
90          % Reset the hold state
91          hold off;
92      end
93  PlotDimensions(figure(1), 'centimeters', [15.747, 14], 12)
94  ChangeInterpreter(figure(1), 'latex')
95
96  % Plot image in pdf format
97  h = figure(1);
98  set(h,'Units','Inches');
99  pos = get(h,'Position');
100 set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
101     'PaperSize', [pos(3), pos(4)])
102 print(h, 'histogram', '-dpdf', '-r0')
103
104 % Calculate the Pearson correlation coefficient between
105 % each pair of features
106 corrMatrix = corr(data);
107
108 % Create a heatmap of the correlation matrix
```

```matlab
109  figure(2);
110  heatmap(featureNames, featureNames, corrMatrix);
111  title('Pearson correlation between features')
112  PlotDimensions(figure(2), 'centimeters', [15.747, 14], 12)
113  ChangeInterpreter(figure(2), 'latex')
114
115  % Plot image in pdf format
116  h = figure(2);
117  set(h,'Units','Inches');
118  pos = get(h,'Position');
119  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
120      'PaperSize', [pos(3), pos(4)])
121  print(h, 'corr1', '-dpdf', '-r0')
122
123  % Calculate the mean and standard deviation of each column
124  meanVals = mean(data, 1);
125  stdDevs = std(data, 0, 1);
126
127  % Perform standard score normalization on each feature
128  standardizedData = (data - meanVals) ./ stdDevs;
129
130  % Calculate the Pearson correlation coefficient between
131  % each pair of standardised features
132  standardizedCorrMatrix = corr(standardizedData);
133
134  % Create a heatmap of the standardised correlation matrix
135  figure(3);
136  heatmap(featureNames, featureNames, standardizedCorrMatrix);
137  title('Pearson correlation between standardised features')
138  PlotDimensions(figure(3), 'centimeters', [15.747, 14], 12)
139  ChangeInterpreter(figure(3), 'latex')
140
141  % Plot image in pdf format
142  h = figure(3);
143  set(h,'Units','Inches');
144  pos = get(h,'Position');
145  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
146      'PaperSize', [pos(3), pos(4)])
147  print(h, 'corr2', '-dpdf', '-r0')
148
149  % Find the minimum and maximum values in each column of the data
150  minVals = min(data);
151  maxVals = max(data);
152
153  % Normalize the data using max-min normalization
154  normalisedData = (data - minVals) ./ (maxVals - minVals);
155
156  % Calculate the Pearson correlation coefficient between
157  % each pair of normalised features
158  minMaxCorrMatrix = corr(normalisedData);
159
160  % Create a heatmap of the min-max correlation matrix
```

```matlab
161  figure(4);
162  heatmap(featureNames, featureNames, minMaxCorrMatrix);
163  title('Pearson correlation between normalised features')
164  PlotDimensions(figure(4), 'centimeters', [15.747, 14], 12)
165  ChangeInterpreter(figure(4), 'latex')
166
167  % Plot image in pdf format
168  h = figure(4);
169  set(h, 'Units', 'Inches');
170  pos = get(h, 'Position');
171  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
172      'PaperSize', [pos(3), pos(4)])
173  print(h, 'corr3', '-dpdf', '-r0')
174
175  %% ================ Part 2: Feature selection
176          ========================
177
177  % Normalize the columns by their mean values
178  %dataNorm = data ./ mean(data, 1);
179
180  % Create the first new feature by adding the normalized
181  % columns 2 and 4
182  %trade = dataNorm(:, 2) + dataNorm(:, 4);
183  trade = data(:, 4);
184  % Create the second new feature by adding the normalized
185  % columns 5, 6, and 9
186  %finance = dataNorm(:, 5) + dataNorm(:, 6) + dataNorm(:, 9);
187  finance = data(:, 6);
188  % Create the third new feature by concatenating the remaining
189          columns
189  %health = dataNorm(:, 1) + dataNorm(:, 3) + dataNorm(:, 7)...
190  %     + dataNorm(:, 8);
191  health = data(:, 3);
192  % Concatenate the new features
193  newFeatures = [trade finance health];
194
195  % Normalise the new features using max-min normalization
196  dataFinal = (newFeatures - min(newFeatures)) ./...
197      (max(newFeatures) - min(newFeatures));
198
199
200  %% ===== Part 3: Selection and execution of clustering algorithms
201          ======
201
202  % Transpose dataFinal for input to k_means function
203  dataFinal = dataFinal';
204
205  % Set number of runs and range of values for m
206  nRuns = 40;
207  mMin = 2;
208  mMax = 10;
209
```

```matlab
210  % Preallocate array to store results
211  jM = zeros(1, mMax - mMin + 1);
212
213  % Loop over values of m
214  for m = mMin:mMax
215      % Initialize temporary minimum value
216      jTempMin = inf;
217
218      % Loop over number of runs
219      for t = 1:nRuns
220          % Generate initial theta values using randDataInit function
221          thetaIni = rand_data_init(dataFinal, m);
222
223          % Run kMeans function and store results
224          [theta, bel, j] = k_means(dataFinal, thetaIni);
225
226          % Update temporary minimum value if necessary
227          if jTempMin > j
228              jTempMin = j;
229          end
230      end
231
232      % Append minimum value to jM array
233      jM(m - mMin + 1) = jTempMin;
234  end
235
236  % Define m values for plot
237  m = mMin:mMax;
238
239  % Create figure and plot jM versus m
240  figure(5), plot(m, jM);
241  xlabel("Number of clusters m");
242  ylabel("Cost function values J");
243  ChangeInterpreter(figure(5), 'latex')
244
245  % Plot image in pdf format
246  h = figure(5);
247  set(h,'Units','Inches');
248  pos = get(h,'Position');
249  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
250      'PaperSize', [pos(3), pos(4)])
251  print(h, 'elbow', '-dpdf', '-r0')
252
253  evaluation = evalclusters(dataFinal', "kmeans",...
254      "CalinskiHarabasz", "KList", 1:10);
255  figure(6), plot(evaluation)
256  ChangeInterpreter(figure(6), 'latex')
257
258  % Plot image in pdf format
259  h = figure(6);
260  set(h,'Units','Inches');
261  pos = get(h,'Position');
```

```matlab
262  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
263      'PaperSize', [pos(3), pos(4)])
264  print(h, 'eval', '-dpdf', '-r0')
265
266  % Set fixed value of m
267  m = 3;
268
269  % Preallocate theta and bel arrays
270  nRuns = 100;
271  theta = zeros(m, size(dataFinal, 1), nRuns);
272  bel = zeros(1, size(dataFinal, 2), nRuns);
273
274  % Initialize temporary minimum value
275  jTempMin = inf;
276
277  % Loop over number of runs
278  for t = 1:nRuns
279      % Generate initial theta values using randDataInit function
280      thetaIni = rand_data_init(dataFinal, m);
281
282      % Run kMeans function and store results
283      [theta(:, :, t), bel(:, :, t), j] = k_means(dataFinal, thetaIni)
          ;
284
285      % Update temporary minimum value and corresponding
286      % outputs if necessary
287      if jTempMin > j
288          jTempMin = j;
289          thetaMin = theta(:, :, t);
290          belMin = bel(:, :, t);
291      end
292  end
293
294  % Plot the clusters
295  figure(7), plot3(dataFinal(1, belMin==1),...
296      dataFinal(2, belMin==1), dataFinal(3, belMin==1),'rx',...
297      dataFinal(1, belMin==2), dataFinal(2, belMin==2),...
298      dataFinal(3, belMin==2),'go', dataFinal(1, belMin==3),...
299      dataFinal(2, belMin==3), dataFinal(3, belMin==3),'bd');
300  hold on
301  plot3(thetaMin(1,:), thetaMin(2,:), thetaMin(3,:), 'k+', 'LineWidth'
      , 2)
302  xlabel("Trade")
303  ylabel("Finance")
304  zlabel("Health")
305  legend("Cluster 1", "Cluster 2", "Cluster 3", "Centroid",...
306      'Location', 'NorthEast')
307  grid on
308  hold off
309
310  ChangeInterpreter(figure(7), 'latex')
311  PlotDimensions(figure(7), 'centimeters', [18, 18], 12)
```

```matlab
312  Plot2LaTeX(figure(7), 'test')
313
314  %% ============= Part 4: Characterisation of clusters
         ================
315
316  cluster1 = countryNames(belMin == 1);
317  cluster2 = countryNames(belMin == 2);
318  cluster3 = countryNames(belMin == 3);
319
320  % Get the lengths of the cell arrays
321  len1 = length(cluster1);
322  len2 = length(cluster2);
323  len3 = length(cluster3);
324
325  % Sort the strings in each cell array in alphabetical order
326  cluster1 = sort(cluster1);
327  cluster2 = sort(cluster2);
328  cluster3 = sort(cluster3);
329
330  % Print the cluster labels
331  fprintf('%20s\t%20s\t%20s\n', 'Cluster 1', 'Cluster 2', 'Cluster 3')
         ;
332  fprintf('%20s\t%20s\t%20s\n', '----------', '----------',...
333      '----------');
334
335  % Print the first 5 strings of each cell array
336  fprintf('%20s\t%20s\t%20s\n', cluster1{1:5}, ...
337      cluster2{1:5}, cluster3{1:5});
338
339  % Print the dots if there are more than 10 strings
340  % in any of the cell arrays
341  if len1 > 10 || len2 > 10 || len3 > 10
342      fprintf('%20s\t%20s\t%20s\n', '.', '.', '.');
343      fprintf('%20s\t%20s\t%20s\n', '.', '.', '.');
344      fprintf('%20s\t%20s\t%20s\n', '.', '.', '.');
345  end
346
347  % Print the last 5 strings of each cell array
348  fprintf('%20s\t%20s\t%20s\n', cluster1{max(1, len1-4):len1},...
349      cluster2{max(1, len2-4):len2}, cluster3{max(1, len3-4):len3});
350
351  clusters = ["Cluster 1", "Cluster 2", "Cluster 3"];
352  clusterData1 = data(belMin' == 1, :);
353  clusterData2 = data(belMin' == 2, :);
354  clusterData3 = data(belMin' == 3, :);
355
356  % Concatenate row vectors and append data to grouping variables
357  childMortality = [clusterData1(:, 1)' clusterData2(:, 1)'...
358      clusterData3(:, 1)'];
359  income = [clusterData1(:, 5)' clusterData2(:, 5)'...
360      clusterData3(:, 5)'];
361  grp = [zeros(1, length(clusterData1(:, 1)')),...
```

```
362        ones(1, length(clusterData2(:, 1)')),...
363        2 * ones(1, length(clusterData3(:, 1)'))];
364
365  % Boxplot of child mortality for all clusters
366  figure(8), boxplot(childMortality, grp, 'Labels', clusters);
367  ylabel("child mortality")
368  ChangeInterpreter(figure(8), 'latex')
369
370  % Plot image in pdf format
371  h = figure(8);
372  set(h,'Units','Inches');
373  pos = get(h,'Position');
374  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
375        'PaperSize', [pos(3), pos(4)])
376  print(h, 'childmortbox', '-dpdf', '-r0')
377
378  % Boxplot of income for all clusters
379  figure(9), boxplot(income, grp, 'Labels', clusters);
380  ylabel("income")
381  ChangeInterpreter(figure(9), 'latex')
382
383  % Plot image in pdf format
384  h = figure(9);
385  set(h,'Units','Inches');
386  pos = get(h,'Position');
387  set(h, 'PaperPositionMode', 'Auto', 'PaperUnits', 'Inches',...
388        'PaperSize', [pos(3), pos(4)])
389  print(h, 'healthbox', '-dpdf', '-r0')
```

Listing 2: MATLAB® function for changing the interpreter of all objects within a figure.

```
1   function ChangeInterpreter(h, Interpreter)
2   % ChangeInterpreter() changes the interpreter of figure h.
3
4       % Find all string type objects
5       TexObj = findall(h, 'Type', 'Text');
6       LegObj = findall(h, 'Type', 'Legend');
7       AxeObj = findall(h, 'Type', 'Axes');
8       ColObj = findall(h, 'Type', 'Colorbar');
9
10      Obj = [TexObj; LegObj]; % Tex and Legend opbjects can be treated
              similarly
11      n_Obj = length(Obj);
12      for i = 1:n_Obj
13          Obj(i).Interpreter = Interpreter;
14      end
15
16      Obj = [AxeObj; ColObj]; % Axes and ColorBar objects can be
              treated similarly
17      n_Obj = length(Obj);
18      for i = 1:n_Obj
19          Obj(i).TickLabelInterpreter = Interpreter;
```

```
20        end
21  end
```

Listing 3: MATLAB® function for configuring a figure's appearance options.

```matlab
function PlotDimensions(h, Units, Plotsize, Fontsize)
% PlotDimensions() changes the string units, the fontsize
% and the unit size of the figure h.

    h.Units = Units; % measurement units
    h.Position(2) = (h.Position(2) - 8.5); % bottom-left corner of
        plot
    h.Position((3:4)) = Plotsize; % usually [15.747, 9]
    set(findall(h, '-property', 'FontSize'), 'FontSize', Fontsize);
        % fontsize
end
```