# Correspondence Analysis

Thomas Bosch
Technical University of Munich, Germany

22.01.2009

**Abstract**

*Correspondence Analysis*
The idea behind correspondence analysis is the optimal representation of a contingency table in low-dimensional space for easy interpretation of any dependency between rows and columns (Bendixen, 1996).

*Sex and Color*
There is no significant dependency between the variables "sex" and "color".

*Sex and Brand*
There is no significant dependency between the variables "sex" and "brand".

*Age and Color*
There is no significant dependency between the variables "age" and "color".

*Age and Brand*
There is observed significant dependency between the variables "age" and "brand".

*Keywords*: Correspondence Analysis; Contingency Table; Asymmetric Plot

## 1.  Introduction

The purpose of this research paper is to apply correspondence analysis to the underlying data set in order to examine how two different demographics sex and age relate to the color and the brand of cars.

The questions to be answered are the following:

1. Are there any interesting relationships?
2. How strong are these relationships?
3. How does this information can be used for marketing purposes by the related companies?

## 2.  Description of the Data

The data set refers to a part of a marketing survey. The respondents were asked with respect to their car.

### 2.1. Individuals

The 251 individuals were interviewed persons in the context of a marketing survey.

### 2.2. Variables

Actually, the data contains only four variables, which relate to

1. sex,
2. age,
3. color (of the car) and
4. brand (of the car).

The possible values of each variable will be described.

*sex*
female, male

*age*
18-25, 25-35, 35-40, 40-45, 45-50, 50+

*color*
white, black, yellow, red, green, blue, grey

*brand*
A, B, C, D, E, F, G

In order to avoid confidentiality issues by the data provider, the cars are referred to as A to G.

## 3.  Correspondence Analysis

The correspondence analysis will be described in general as a first step. Then the formal procedure will be shown. Finally, the methodology will be executed with the given data.

### 3.1. Description

The idea behind correspondence analysis is the optimal representation of a contingency table in low-dimensional space for easy interpretation of any dependency between rows and columns (Bendixen, 1996).

*3.2. Methodology*

*Relationships*
The relationships between

- sex and color,
- sex and brand,
- age and color and
- age and brand

will be examined.

*Statistic Software R*
The data will be imported with the statistic software R in its current version 2.7.2. All the R statements used to get necessary results can be read in the appendix of this paper.

*Data Import to R*
The data is given in a SPSS format. For this reason you have to install an additional package in R. The package "foreign" provides functions for the data import from SPSS to R. Before using this package, it has to be activated. Then the data can be read.

*Contingency Tables*
A contingency table shows the frequency of congruence of the values of two dimensions among the sample (Bendixen, 1996).

The contingency tables for each relationship will be calculated in excel. The variables from the data set will be exported from the tool R to an excel file. The values of these variables will be imported to another excel file, which will be utilized for the following calculations. In this excel file, the contingency tables will be generated using visual basic for applications. The appropriate code listing can be read in the appendix.

*XLSTAT*
In order to execute the correspondence analysis, the statistical excel tool XLSTAT was used. You can download an evaluation version from the website of this tool. First of all, you have to import the contingency table to another excel file. After the installation of XLSTAT, you have the choice between different toolbars in excel. Click on the "Analyzing Data" button of the "Analyzing Data" toolbar. Once you have clicked on the button, the correspondence analysis dialog box appears. Then select the data on the Excel sheet. If your data are in a pivot table format, select the contingency table format. If your data are in an observations/variables format, select the corresponding option. As the names of the categories of the contingency table were included, the "Labels included" option was selected as well. The author of this research paper selected the "Range" option for the output. The $J$4 cell was selected as the upper left corner of the results report. You can also choose to write the results in a separate sheet or workbook. On the charts tab, the author selected to display the asymmetric rows and columns plot. In order to get all the necessary information, "additional data" in the options tab was selected. The results are displayed once the user has selected and validated the axes on which the plots need to be displayed.

*Significance of Dependencies between rows and columns*
"The first step in the interpretation of correspondence analysis is to establish weather there is a significant dependency between the rows and columns. There are two approaches to establish significance [: the trace and the chi-square test of independence]" (Bendixen, 1996, p. 25).

*Trace*
The trace is the sum of the eigenvalues of the axes. "The square root of the trace may be interpreted as a correlation coefficient between the rows and columns (Bendixen, 1996, p. 25). Values greater than 0.2 indicate significant dependency. "This is a rough and ready approximation and a more thorough approach is to calculate the chi-square statistic […]" (Bendixen, 1996, p. 25).

*Chi-Square Test of Independence*
"This statistical test is used to determine […] whether there is a statistically significant dependence between the rows and the columns [of the contingency table]" (Bendixen, 1996, p. 17). There are two hypotheses. The first one stands for the fact that the rows and the columns of the contingency table are not dependent and the second that the rows and the columns are dependent. If the calculated p-value is smaller than the level of significance alpha, the null-hypothesis must be rejected and the alternative hypothesis must be accepted.

*Dimensionality of the solution*
You have to determine the number of dimensions used in the solution. You have to examine the eigenvalue report (Bendixen, 1996). "The ratio of the eigenvalue of any axis to the trace represents the proportion of the total "inertia" […] explained by that axis" (Bendixen, 1996, p. 26). In terms of the rows, the average axis should account for 100 / (< number of rows > - 1) of the inertia. In terms of the columns, the average axis should account for 100 / (< number of columns > - 1) of the inertia (Bendixen, 1996). "[…] Any axis contributing more than the maximum of these two percentages should be regarded as significant and included in the solution" (Bendixen, 1996, p. 26). "[…] A higher number of dimensions may be used but the additional dimensions are unlikely to contribute significantly to the interpretation of [the] nature of the dependency between [the] rows and [the] columns" (Bendixen, 1996, p.26). The contribution of the selected axes to the inertia is called the retention of the solution (Bendixen, 1996). "The higher the retention, the more subtlety in the original data is retained in the low-dimensional solution" (Bendixen, 1996, p. 26).

*Interpreting the axes*
In this research paper, just asymmetric plots and not symmetric plots are used, because symmetric plots can easily lead to misunderstanding. "The apices of either the rows (or the columns) are plotted from the standard coordinates and the profiles of the columns (or the rows) are plotted from the principle coordinates (Bendixen, 1996, p. 27).

The next step "[…] is to interpret the axis in terms of the rows (or the columns) and [to] plot [..] the column points (or [the] row points) in the space of the labeled axes (Bendixen, 1996, p. 27). You have "[…] to decide whether to interpret the axes in terms of rows or columns (Bendixen, 1996, p. 27).

In the next step, "the axes are interpreted by way of the contribution […] [of each row (or column) to] the total inertia accounted for by the axis" (Bendixen, 1996, p. 28). "Any contribution greater than 100 / [..] [<number of rows or columns>] would represent significance greater than what would be expected in the case of a purely random distribution of [..] [rows or columns] over the axes" (Bendixen, 1996, p. 28).

The principle coordinates can now be plotted with the axes labeled as determined.

*Quality of Representation*
"[..] Not all of the [..] [rows] or [..] [columns] are equally well represented. Determining the quality of representation of a particular row or column provides additional richness to the interpretation of the relationships in the contingency table" (Bendixen, 1996, p. 30). "[…]

The quality of representation is easily calculated from the […] [squared cosine values] given in the output" (Bendixen, 1996, p. 30). "The […] [squared cosine value] presented for any [..] [row (or column)] measures the degree of association between that [..] [row (or column)] and a particular axis" (Bendixen, 1996, p. 30).

"The quality of representation of a row or column in n dimensions is simply the sum of the […] [squared cosine values] of that row or column over the n dimensions" (Bendixen, 1996, p. 31). If the quality of representation is high enough, the row or the column is well presented in low-dimensional space. If the quality of representation is low, a higher dimensional solution is probably necessary to understand the relationship between the rows and the columns (Bendixen, 1996).

*Marketing*
In this section the question will be answered how the collected information can be used for marketing purposes by the related companies.

*Conclusion*
In conclusion, the most important findings will be summarized.

## 3.3. Relationship between Sex and Color

The relationship between the variables "sex" and "color" will be executed.

### 3.3.1.    Contingency Table

The contingency table of the association between the variables "sex" and "color" is shown in the following table.

Table 1
*Contingency Table of Sex and Color*

|        | white | black | red | green | grey |
|--------|-------|-------|-----|-------|------|
| female | 23    | 21    | 25  | 17    | 17   |
| male   | 26    | 14    | 12  | 16    | 12   |

Since row and column sums do not have to be zero, the columns "yellow" and "blue" were deleted from the contingency table.

### 3.3.2.    Significance of Dependencies between rows and columns

*Trace*
The trace is 0.02 and the square root of the trace is 0.15. Since this value is less than 0.2, the rows ant the columns are not dependent.

*Chi-Square Test of Independence*
Since the calculated p-value 0.377 is higher than the level of significance alpha 0.05, the null-hypothesis must be accepted. There is observed no significant dependency between the rows and the columns of the contingency table.

Both the square root of the trace and the chi-square test of independence show no significant dependency.

### 3.3.3.    Conclusion

Since there is no significant dependency between the two variables "sex" and "color", the further steps of a correspondence analysis will not be executed.

## 3.4. Relationship between Sex and Brand

The relationship between the variables "sex" and "brand" will be executed.

### 3.4.1. Contingency Table

The contingency table of the association between the variables "sex" and "brand" is shown in the following table.

Table 2
*Contingency Table of Sex and Brand*

|  | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| female | 29 | 15 | 14 | 23 | 12 | 29 | 17 |
| male | 23 | 16 | 12 | 8 | 18 | 22 | 13 |

### 3.4.2. Significance of Dependencies between rows and columns

*Trace*
The trace is 0.032 and the square root of the trace is 0.179. Since this value is less than 0.2, the rows ant the columns are not dependent.

*Chi-Square Test of Independence*
Since the calculated p-value 0.237 is higher than the level of significance alpha 0.05, the null-hypothesis must be accepted. There is observed no significant dependency between the rows and the columns of the contingency table.

Both the square root of the trace and the chi-square test of independence show no significant dependency.

### 3.4.3. Conclusion

Since there is no significant dependency between the two variables "sex" and "brand", the further steps of a correspondence analysis will not be executed.

### 3.5. Relationship between Age and Color

The relationship between the variables "age" and "color" will be executed.

### 3.5.1. Contingency Table

The contingency table of the association between the variables "age" and "color" is shown in the following table.

Table 3
*Contingency Table of Age and Color*

|  | white | black | red | green | grey |
|---|---|---|---|---|---|
| 18-25 | 19 | 12 | 17 | 10 | 7 |
| 25-35 | 16 | 18 | 16 | 17 | 16 |
| 35-40 | 8 | 2 | 2 | 5 | 5 |
| 40-45 | 5 | 2 | 1 | 0 | 1 |
| 45-50 | 1 | 1 | 0 | 0 | 0 |
| 50+ | 0 | 0 | 1 | 1 | 0 |

Since row and column sums do not have to be zero, the columns "yellow" and "blue" were deleted from the contingency table.

### 3.5.2. Significance of Dependencies between rows and columns

*Trace*
The trace is 0.11 and the square root of the trace is 0.33 indicating significant dependency.

*Chi-Square Test of Independence*
Since the calculated p-value 0.432 is higher than the level of significance alpha 0.05, the null-hypothesis must be accepted. There is observed no significant dependency between the rows and the columns of the contingency table.

According to the square root of the trace, there is a significant dependency between the rows and the columns, and corresponding to the chi-square test of independence, there is no significant dependency observed between the rows and the columns.

### 3.5.3. Conclusion

The chi-square test of independence is more exact in the determination of the dependency between the rows and the columns of the contingency table. Because of this, there is no significant dependency between the two variables "age" and "color", the further steps of a correspondence analysis will not be executed.

### 3.6. Relationship between Age and Brand

The relationship between the variables "age" and "brand" will be executed.

### 3.6.1. Contingency Table

The contingency table of the association between the variables "age" and "brand" is shown in the following table.

Table 4
*Contingency Table of Age and Brand*

|  | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 18-25 | 27 | 10 | 11 | 9 | 12 | 21 | 6 |
| 25-35 | 13 | 19 | 12 | 15 | 14 | 20 | 16 |
| 35-40 | 10 | 2 | 1 | 2 | 3 | 5 | 5 |
| 40-45 | 2 | 0 | 2 | 5 | 1 | 1 | 2 |
| 45-50 | 0 | 0 | 0 | 0 | 0 | 2 | 1 |
| 50+ | 0 | 0 | 0 | 0 | 0 | 2 | 0 |

### 3.6.2. Significance of Dependencies between rows and columns

*Trace*
The trace is 0.18 and the square root of the trace is 0.42 indicating significant dependency.

*Chi-Square Test of Independence*

Since the calculated p-value 0.0037 is smaller than the level of significance alpha 0.05, the null-hypothesis must be rejected and the alternative hypothesis must be accepted. There is observed significant dependency between the rows and the columns of the contingency table.

Both the square root of the trace and the chi-square test of independence show significant dependency.

### 3.6.3. Dimensionality of the Solution

In terms of the rows, the average axis should account for 100 / (6 - 1) = 20 percent of the inertia, because the number of rows is 6. In terms of the columns, the average axis should account for 100 / (7 - 1) = 16.6 percent of the inertia, because the number of columns is 7. Any axis contributing more than the maximum of these two percentages should be regarded as significant and included in the solution. The third axis accounts only 17.14 percent of the inertia. Because of this, a two-dimensional solution should be used. The first and the second axis account for 39.5 percent and 30.31 percent of the inertia that is a cumulative total of 69.81 percent also called the retention.

### 3.6.4. Interpreting the Axes

*Asymmetric Plots*

In the following figures you can see the asymmetric plots. The first one shows the interpretation of the axes in terms of the rows and the second one represents the interpretation of the axes in terms of the columns.
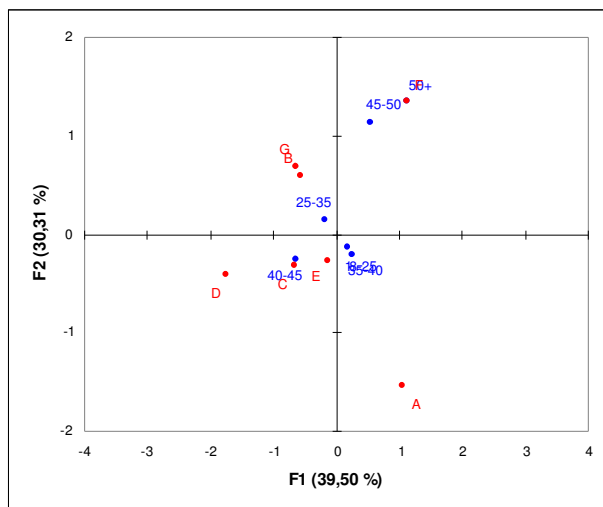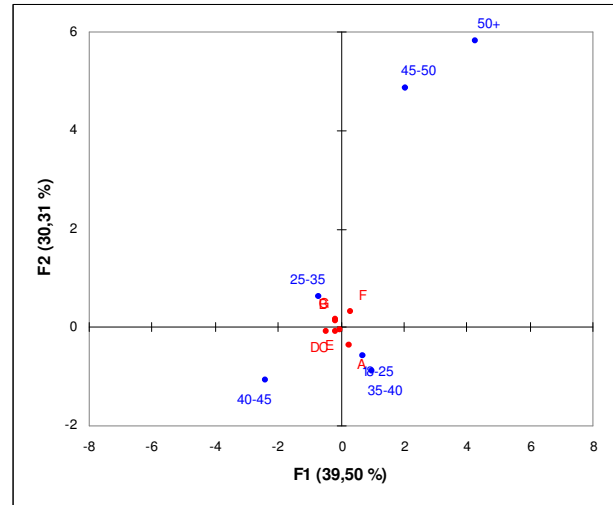


Figure 2
*Asymmetric Plot of Ages in Brand Space*

In order to determine the brands of specific age clusters, it seems to be most appropriate to interpret the axes in terms of the rows that is brands in age space. The individual car manufacturers have the possibility to investigate the target groups in terms of the age of their brand.

*Determination of the axes*

Any contribution greater than 100 / 6 = 16.6 percent would represent significance greater than what would be expected in the case of a purely random distribution of ages over the axes, because the number of ages is 6.

*Determination of the first axis*

Examining the detailed report for the rows, the ages 18-25, 25-35 and 40-45 meet this criterion and determine the first axis. The ages 18-25 and 40-45 have positive coordinates and the ages 25-35 have negative coordinates.

*Determination of the second axis*

The ages 25-35, 45-50 and 50+ meet the criterion, too, and determine the second axis. The ages 25-35, 45-50 and 50+ have positive coordinates. There are no age groups having negative coordinates.

*Asymmetric Plot of Brands in Age Space*

The resulting asymmetric plot is displayed in the following figure.



Figure 1
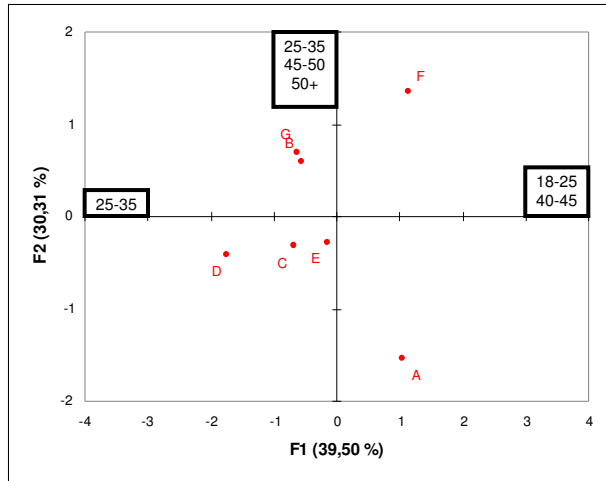*Asymmetric Plot of Brands in Age Space*

Figure 3
*Asymmetric Plot of Brands in Age Space*

You can see that cars with the brand name B and G are mostly driven by persons whose age is between 25 and 35 or between 45 and 50 or over 50. Brands C, D and E are in general driven by humans with the age between 25 and 35. The brand F is in the most cases driven by persons of all age clusters except 35-40. Humans with the age of either between 18 and 25 or between 40 and 45 drive the car of the brand A.

*3.6.5. Quality of Representation*

The quality of representation of every row will be determined.

*Brand A*
The squared cosine value between brand A and the first and second axis is 0.373 and 0.612. This implies that brand A is strongly associated with the second axis but only weakly associated with the first axis. Brand A is strongly associated with drivers, which are not part of the three clusters 25-35, 45-50 and 50+. Brand A is weakly associated with drivers between 10 and 25 or between 40 and 45.

The quality of representation of brand A is 0.373 + 0.612 = 0.985. Brand A is well presented in the two dimensions.

*Brand B*
The squared cosine value between brand B and the first and second axes is 0.129 and 0.113. This implies that brand B is weakly associated with the first and the second axis. Brand B is weakly associated with drivers between 25 and 35, between 45 and 50 and over 50.

The quality of representation of brand B is 0.129 + 0.113 = 0.242. Brand B is not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between brand B and the rows.

*Brand C*
The squared cosine value between brand C and the first and second axes is 0.377 and 0.062. This implies that brand C is weakly associated with the first and the second axis. Brand C is weakly associated with drivers between 25 and 35 and drivers, which are not part of the age groups 25-35, 45-50 and over 50.

The quality of representation of brand C is 0.377 + 0.062 = 0.439. Brand C is not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between brand C and the rows.

*Brand D*
The squared cosine value between brand D and the first and second axes is 0.720 and 0.030. This implies that brand D is strongly associated with the first axis but only weakly associated with the second axis. Brand D is strongly associated with drivers, which are between 25 and 35 years old. Brand D is weakly associated with drivers which are not part of the age groups 25-35, 45-50 and over 50.

The quality of representation of brand D is 0.720 + 0.030 = 0.75. Brand D is well presented in the two dimensions.

*Brand E*
The squared cosine value between brand E and the first and second axes is 0.043 and 0.133. This implies that brand E weakly associated with the first and the second axis. Brand E is weakly associated with drivers, which are between 25 and 35 years old. Brand E is also weakly associated with drivers which are not part of the age clusters 25-35, 45-50 and over 50.

The quality of representation of brand E is 0.043 + 0.133 = 0.176. Brand E is not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between brand E and the rows.

*Brand F*
The squared cosine value between brand F and the first and second axes is 0.430 and 0.472. This implies that brand F is weakly associated with the first and the second axis. Brand F is weakly associated with drivers between 18 and 25 or between 40 and 45. Brand F is also weakly associated with drivers which are part of the age clusters 25-35, 45-50 and over 50.

The quality of representation of brand F is 0.430 + 0.472 = 0.902. Brand F is well presented in the two dimensions.

*Brand G*
The squared cosine value between brand G and the first and second axis is 0.155 and 0.139. This implies that brand G is weakly associated with the first and the second axis. Brand G is weakly associated with drivers between 25 and 35. Brand G is also weakly associated with drivers which are not part of the age clusters 25-35, 45-50 and over 50.

The quality of representation of brand G is 0.155 + 0.139 = 0.294. Brand G is not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between brand G and the rows.

*Good Quality of Representation*
The brands A, D and F are well presented in the two dimensions.

*Bad Quality of Representation*
The brands B, C, E and G are not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between these brands and the age clusters.

*3.6.6. Marketing*

The marketers of the car manufacturers can determine which brands are driven by persons assigned to specific age clusters. The individual car manufacturers have the possibility to investigate the target groups in terms of the age of their brand. The marketers see, if the current target group of the appropriate age segment matches with the target group they want. If this is not the case, the manufactures of the brands must change their marketing mix or their strategy relating to the specific brand. It is very important for the marketing executives to know if their strategy is implemented in the way they defined it in the past. So they have the chance to react earlier.

*3.6.7. Conclusion*

The relationship between the variables "age" and "brand" was executed. The most important findings will be summarized.

*Significance of Dependencies between rows and columns*
Both the square root of the trace and the chi-square test of independence showed significant dependency between the rows and the columns of the underlying contingency table.

*Dimensionality of the solution*
The third axis accounts only 17.14 percent of the inertia. Because of this, a two-dimensional solution should be used. The first and the second axis account for 39.5 percent and 30.31 percent of the inertia that is a cumulative total of 69.81 percent also called the retention.

*Interpreting the axes*
The axes are interpreted in terms of the rows that is brands in age space. You can see the axes labeled with the age clusters in the following asymmetric plot.
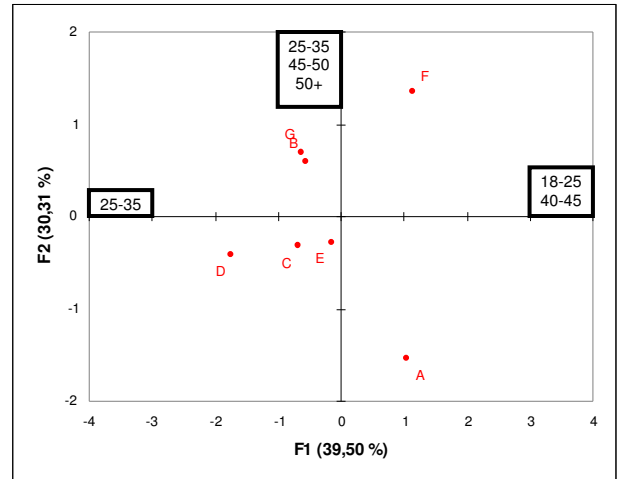


Figure 4
*Asymmetric Plot of Brands in Age Space*

*Quality of Representation*
The brands A, D and F are well presented in the two dimensions.

The brands B, C, E and G are not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between these brands and the age clusters.

*Marketing*
The car manufacturers have the possibility to determine the target groups of their brand in terms of the age.

**4. Conclusion**

*Correspondence Analysis*
The idea behind correspondence analysis is the optimal representation of a contingency table in low-dimensional space for easy interpretation of any dependency between rows and columns (Bendixen, 1996).

*Relationships*
The relationships between

- sex and color,
- sex and brand,
- age and color and
- age and brand

were examined.

*Sex and Color*
There is no significant dependency between the variables "sex" and "color".

*Sex and Brand*
There is no significant dependency between the variables "sex" and "brand".

*Age and Color*
There is no significant dependency between the variables "age" and "color".

*Age and Brand*
There is observed significant dependency between the variables "age" and "brand".

The axes are interpreted in terms of the rows that is brands in age space. You can see the axes labeled with the age clusters in the following asymmetric plot.
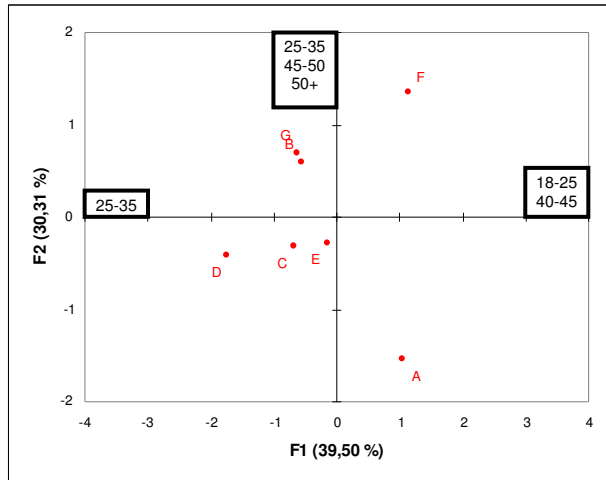


Figure 5
*Asymmetric Plot of Brands in Age Space*

The brands A, D and F are well presented in the two dimensions.

The brands B, C, E and G are not well presented in the two dimensions. A higher dimensional solution is probably necessary to understand the relationship between these brands and the age clusters.

The car manufacturers have the possibility to determine the target groups of their brand in terms of the age.

## Acknowledgment

## References

Bendixen M. (1996). *A Practical Guide to the Use of Correspondence Analysis in Marketing Research.* South Africa: Marketing Research On-Line.

## Appendix

### R Code Listing

```
#----- install and activate package necessary to read from spss -----
install.packages("foreign", dependencies = TRUE)
library(foreign)
#-----

#----- read data from spss -----
data      <-      read.spss("D:/Studium/Erasmus/Courses/Multivariate
Techniques/project 3/project3.sav")
data
#-----

#----- export data to excel -----
write.table(data$sex, file = "excel import.csv")
write.table(data$age, file = "excel import.csv")
write.table(data$color, file = "excel import.csv")
write.table(data$brand, file = "excel import.csv")
#-----
```

### VBA Code Listing

```
Option Explicit

Const numberIndividuals = 251
Const firstRowMatrix As Integer = 2
Const lastRowMatrix As Integer = numberIndividuals + 1

Sub calculate()
    Dim r As Integer
    Dim rm As Integer
    Dim c As Integer
    Dim count As Integer
    Dim firstRow As Integer
    Dim lastRow As Integer
    Dim firstCol As Integer
    Dim lastCol As Integer

    firstCol = 7
    lastCol = 13

    '----- sex - color -----
    firstRow = 3
    lastRow = 4

    For r = firstRow To lastRow
        For c = firstCol To lastCol
            count = 0
            For rm = firstRowMatrix To lastRowMatrix
                If ((Cells(rm, 1).Value = Cells(r, firstCol - 1).Value) And
(Cells(rm, 3).Value = Cells(firstRow - 1, c).Value)) Then
                    count = count + 1
                End If
            Next
            Cells(r, c).Value = count
        Next
    Next
    '-----

    '----- sex - brand -----
    firstRow = 8
```

```
    lastRow = 9

    For r = firstRow To lastRow
        For c = firstCol To lastCol
            count = 0
            For rm = firstRowMatrix To lastRowMatrix
                If ((Cells(rm, 1).Value = Cells(r, firstCol - 1).Value) And
(Cells(rm, 4).Value = Cells(firstRow - 1, c).Value)) Then
                    count = count + 1
                End If
            Next
            Cells(r, c).Value = count
        Next
    Next
    '-----

    '----- age - color -----
    firstRow = 13
    lastRow = 18

    For r = firstRow To lastRow
        For c = firstCol To lastCol
            count = 0
            For rm = firstRowMatrix To lastRowMatrix
                If ((Cells(rm, 2).Value = Cells(r, firstCol - 1).Value) And
(Cells(rm, 3).Value = Cells(firstRow - 1, c).Value)) Then
                    count = count + 1
                End If
            Next
            Cells(r, c).Value = count
        Next
    Next
    '-----

    '----- age - brand -----
    firstRow = 22
    lastRow = 27

    For r = firstRow To lastRow
        For c = firstCol To lastCol
            count = 0
            For rm = firstRowMatrix To lastRowMatrix
                If ((Cells(rm, 2).Value = Cells(r, firstCol - 1).Value) And
(Cells(rm, 4).Value = Cells(firstRow - 1, c).Value)) Then
                    count = count + 1
                End If
            Next
            Cells(r, c).Value = count
        Next
    Next
    '-----
End Sub
```