# Correspondence Analysis in Categorical Data

Sebastián David Reyes Santamaría

March 19, 2025

## 1 Principal Component Analysis

We begin with the principal component analysis for each dataset (men and women). To do this, we perform the decomposition of the covariance matrix $S$ into its eigenvectors and eigenvalues:



**Figure 1:** Decomposition for men



**Figure 2:** Decomposition for women

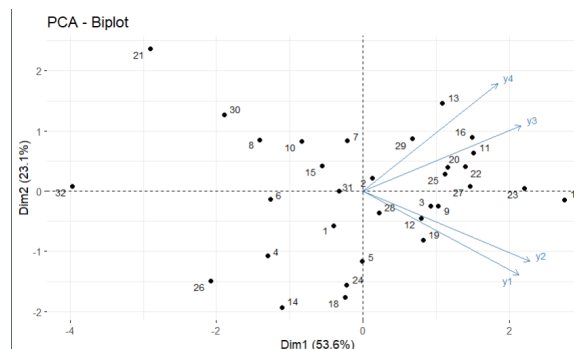The eigenvalues represent the amount of variance explained by each principal component. In this case, the values are higher for the women's dataset, suggesting that the variables might be more correlated.

Visualizing the PCA analysis along with the observations of each dataset, we obtain:
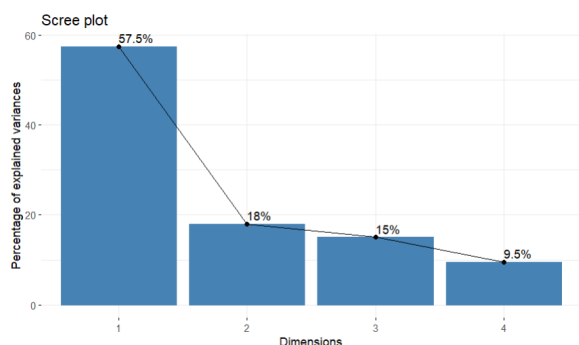
**Figure 3:** PCA for men
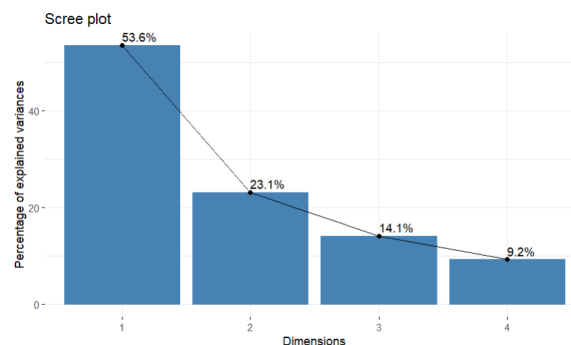


**Figure 4:** PCA for women

Both analyses capture approximately 75% of the total variability of the data in their first two dimensions. According to the results, the variables seem to form two main groups: $(y_3, y_4)$ and $(y_1, y_2)$. In the first component, the explained variance is higher for men (57.50%) than for women (53.56%), whereas in the second component, the explained variance is higher for women (23.10%) than for men (17.97%).

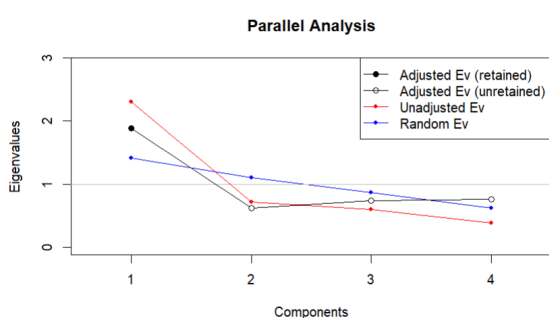**Selection of the number of principal components**

To determine the number of principal components to retain, various criteria were used, such as the rule of eigenvalues greater than 1, the elbow method, and the parallel method. As a result, it was decided to retain only one principal component in both cases.
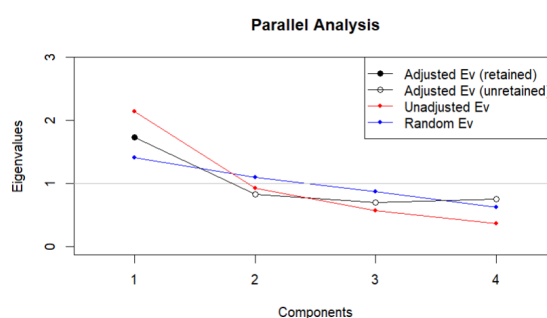
**Figure 5:** Elbow method for men



**Figure 6:** Elbow method for women



**Figure 7:** Parallel method for men



**Figure 8:** Parallel method for women

Under the first criterion of considering components with eigenvalues greater than 1, there is some doubt about the eigenvalue of the second component in the PCA for women, which takes a value of 0.924 (the closest to 1). However, after analyzing with the other two criteria, it is still suggested to select only one component. With the elbow method, a clear break is observed for men in the first component, while for women, it is less pronounced but still follows the expected pattern. Also, Horn's parallel analysis suggests retaining only one component in both cases, as the random eigenvalue line crosses the other lines after the first component.
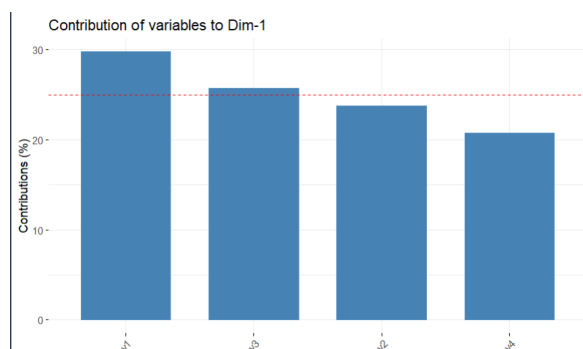
**Interpretation of the Principal Component**

The first component is characterized by a clear opposition between two groups of variables:

- y3 and y4 are positively correlated with each other and have a positive correlation with Dim1.

- y1 and y2 are positively correlated with each other and have a negative correlation with Dim1.
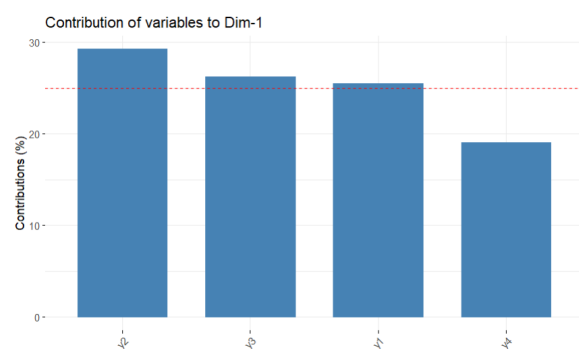
Individuals on the positive side of Dim1 (right) tend to have high values in y3 and y4, and low values in y1 and y2. Conversely, individuals on the negative side of Dim1 (left) tend to have high values in y1 and y2, and low values in y3 and y4.

**Contribution of variables to the principal component**



**Figure 9:** Contribution of each variable in the first principal component for men



**Figure 10:** Contribution of each variable in the first principal component for women

The dominant variable changes: y1 in the first dataset vs. y2 in the second. y4 consistently shows the lowest contribution, although it remains relevant, suggesting that the first principal component captures a robust pattern in the data.

## 2 Correspondence Analysis

For correspondence analysis, contingency tables are first entered into R:



**Figure 11:** Contingency table for womens



**Figure 12:** Contingency table for men

### 2.1 Correspondence analysis for women

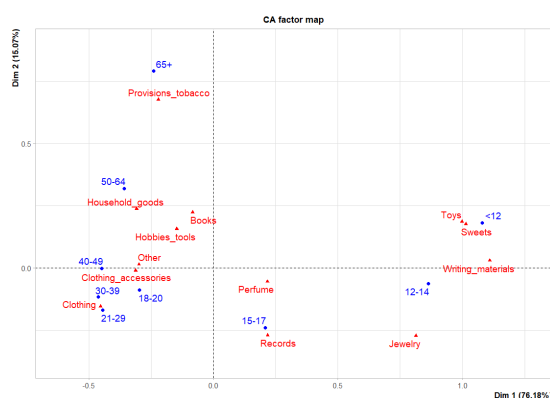**How many dimensions to consider:**

The extracted dimensions explain the variance in the data. Dimension 1 explains 76.18% of the variance, followed by Dimension 2 with 15.07%. Together, the first two dimensions explain 91.25% of the total variance in the data, and since it is recommended to select a sufficient number of dimensions that capture a significant amount of variance in the

4

data, without including dimensions that contribute very little explained variance, we will select these first 2 dimensions for the analysis.

### Main findings

Starting the analysis, we see that the age categories (<12) and (12-14) have the highest contributions to Dimension 1, suggesting that these categories are strongly associated with the (Items) categories in this dimension. The ages (+65) and (50-64) have considerable contributions to Dimension 2, indicating a significant association with different (Items) in this dimension.



**Figure 13:** Correspondence analysis for women

When interpreting Dimension 1, it seems to be more related to age group differences, especially with the older age group (65+), suggesting that people aged 65 or older have distinctive characteristics or preferences compared to other age groups. The younger categories, (15-17) and (18-20), have lower contributions in this dimension. Dimension 2 is more influenced by the Item (Provisions tobacco), in the ages (65+) and (50-64) showing a higher association.

We can see relationships between the categories, for example, that the recorded thefts for ages (<12) and (12-14) tend to be from the Items (Toys, Sweets, Writing materials), for the age (15-17) the theft of the item (Records) is more common, ages (21-29) and (30-39) with the Item (Clothing), age (50-64) with (Household goods) and finally age (65+) with (Provisions tobacco).

## 2.2   Correspondence analysis for men
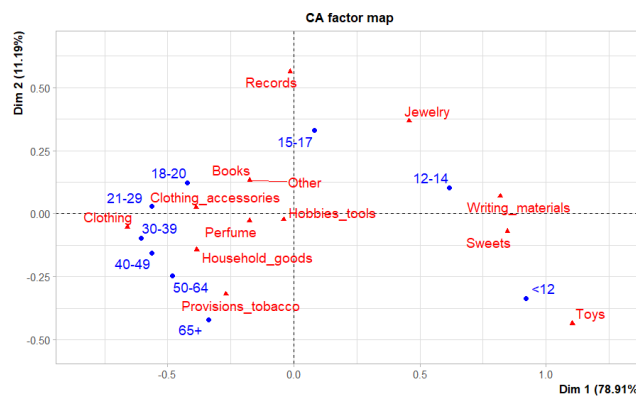
**How many dimensions to consider:**

Similar to the previous analysis, Dimension 1 explains 77.10% of the variance, followed

by Dimension 2 with 11.45%. Together, the first two dimensions explain 88.55% of the total variance in the data. We will select these first 2 dimensions for the analysis.

### Main findings

The age categories (<12, 12-14, 20-29) and the Item (Toys) have the highest contributions to Dimension 1. On the other hand, the age (15-17) and the Item (Records) have a considerable contribution to Dimension 2.
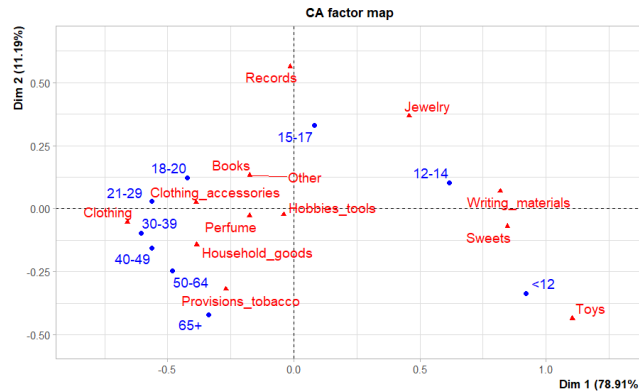


**Figure 14:** Correspondence analysis for men

The younger ages (<12, 12-14) appear to be more related to the Items (Toys, Sweets, Writing materials). On the other hand, the ages (21-29, 30-39) seem to be more associated with the Item (Clothing). Moreover, the older ages, particularly (50-64, 65+) show a strong association with the Items (Household goods and Provisions tobacco), similar to the case of women.

## 2.3  Joint correspondence analysis

Although the positions and contributions of the categories to each dimension vary, in general, a similar association between the different categories is maintained, as seen in the graph. Likewise, the first two dimensions explain 90.1% of the variance in the data.

**Figure 15:** Joint correspondence analysis

The strongest associations remain, for example, for age (<12) being closer to the Item (Toys), age (12-14) to the Items (Sweets and Writing materials) or age (65+) to (Provisions tobacco), in general, maintaining the relationships given in the two separate correspondence analyses for men and women.

## 3   Code

Click here to go to the GitHub repository with the code

File name: CA_analysis_gender_age_groups.R

## References

Heijden, P. G. M., Falguerolles, A., & De Leeuw, J. (2018). A Combined Approach to Contingency Table Analysis Using Correspondence Analysis and Loglinear Analysis. *Journal of the Royal Statistical Society Series C: Applied Statistics*, *38*(2), 249–273. https://doi.org/10.2307/2348058