

Correspondence analysis: introduction

Correspondence Analysis (CA) is a multivariate statistical technique used to analyze and visualize the relationships between categorical variables in a contingency table. It reduces the dimensionality of the data, representing the associations between rows and columns in a low-dimensional space, typically two dimensions, for an easier interpretation. The Chi-squared distance between rows or columns to highlight associations:

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^m \frac{(n_{ij} - Np_{ij})^2}{Np_{ij}}$$

where n_{ij} is the observed frequency and Np_{ij} is the expected frequency under the independence assumption.

'HairEyeColor' dataset

HairEyeColor: dataset of 592 observations x 3 variables.

Hair: qualitative variable: Black, Brown, Red, Blond

Eye: qualitative variable: Brown, Blue, Hazel, Green

Sex: qualitative variable: Male, Female

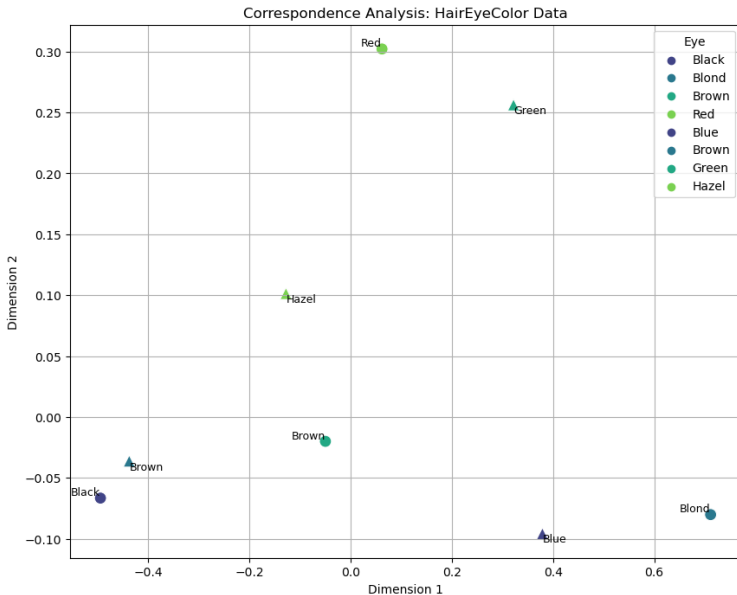
```
1 data = {
2     'Hair': ['Black', 'Brown', 'Red', 'Blond', 'Black', 'Brown', 'Red', 'Blond',
3             'Black', 'Brown', 'Red', 'Blond', 'Black', 'Brown', 'Red', 'Blond'],
4     'Eye': ['Brown', 'Brown', 'Brown', 'Brown', 'Blue', 'Blue', 'Blue', 'Blue',
5            'Hazel', 'Hazel', 'Hazel', 'Hazel', 'Green', 'Green', 'Green', 'Green'],
6     'Freq': [32, 53, 10, 3, 11, 50, 10, 30, 10, 25, 7, 5, 3, 15, 7, 8]
7 }
8
9 hair_eye = pd.DataFrame(data)
10 contingency_table = hair_eye.pivot_table(values='Freq', index='Hair', columns='
    Eye', fill_value=0)
11 contingency_table
12
13 Eye Blue Brown Green Hazel
14 Hair
15 Black 11 32 3 10
16 Blond 30 3 8 5
17 Brown 50 53 15 25
18 Red 10 10 7 7
```

Summary

To perform a correspondence analysis, we can create the function:

```
1 # Correspondence Analysis
2 def correspondence_analysis(table):
3     table = table.values
4     n = np.sum(table)
5     P = table / n
6     r = np.sum(P, axis=1)
7     c = np.sum(P, axis=0)
8     R = np.diag(r)
9     C = np.diag(c)
10    S = P - np.outer(r, c)
11    R_inv_sqrt = np.linalg.inv(np.sqrt(R))
12    C_inv_sqrt = np.linalg.inv(np.sqrt(C))
13    Z = R_inv_sqrt @ S @ C_inv_sqrt
14    U, D, Vt = svd(Z, full_matrices=False)
15    F = R_inv_sqrt @ U @ np.diag(D)
16    G = C_inv_sqrt @ Vt.T @ np.diag(D)
17    return F, G, D
18
19 F, G, D = correspondence_analysis(contingency_table)
20
21 row_coords = pd.DataFrame(F[:, :2], columns=['Dim1', 'Dim2'])
22 row_coords['Hair'] = contingency_table.index
23
24 col_coords = pd.DataFrame(G[:, :2], columns=['Dim1', 'Dim2'])
25 col_coords['Eye'] = contingency_table.columns
```

Plot of factors in main dimentions



Main observations

- Dimension Reduction: The relationships between hair color and eye color in a lower-dimensional space, here the first two on the plot.
- Association Visualization: Points close to each other in the plot indicate a stronger association between the corresponding hair and eye colors. For example, if "Black Hair" and "Brown Eyes" are close together (frequently observed together).
- Dimensional Interpretation: The axes (Dimension 1 and Dimension 2) represent the principal dimensions that capture the most variance in the data.
- Categorical Differentiation: The plot visually differentiates between hair and eye colors using different shapes and colors, making it easy to interpret the correspondence between categories.

References

An Introduction to Applied Multivariate Analysis with R, 2011, B. Everitt, T. Hothorn, Springer, e-ISBN 978-1-4419-9650-3

Python:

<https://www.python.org/>