

# 1 AssociationExplorer: A user-friendly Shiny application 2 for exploring associations and visual patterns

3 Antoine Soetewey<sup>a,b,\*</sup>, Cédric Heuchenne<sup>a,b</sup>, Arnaud Claes<sup>c</sup>, Antonin  
4 Descampe<sup>c</sup>

5 <sup>a</sup>HEC Liège, ULiège, Rue Louvrex 14, 4000 Liège, Belgium

6 <sup>b</sup>The Center for Applied Public Economics (CAPE), UCLouvain Saint-Louis Bruxelles,  
7 Boulevard du Jardin Botanique 43, 1000 Brussels, Belgium

8 <sup>c</sup>Observatory for Research on Media and Journalism (ORM), UCLouvain, Ruelle de la  
9 Lanterne Magique 14, 1348 Louvain-la-Neuve

---

## 10 Abstract

*AssociationExplorer is an interactive R Shiny web application designed to help non-technical users explore statistical associations within multivariate datasets. Aimed particularly at journalists, educators, and engaged citizens, the tool facilitates the discovery and interpretation of meaningful patterns between variables without requiring programming or statistical expertise. Users can upload structured data (e.g., from surveys or open government datasets), select relevant variables, and dynamically visualize relationships via a correlation network and contextual bivariate plots. To illustrate its capabilities, we present a case study based on the European Social Survey (ESS), showcasing how users can investigate links between attitudes, behaviors, and socio-demographic indicators across countries. The app supports a range of association measures adapted to variable types (Pearson's  $r$ , Eta, and Cramer's  $V$ ), ensuring both flexibility and statistical rigor. The visual interface enables users to adjust thresholds for association strength and examine results through interactive graphs and summary tables, making the app particularly well-suited for data storytelling, exploratory research, and public communication. AssociationExplorer demonstrates how open-source statistical tooling can enhance transparency, accessibility, and insight in the interpretation of complex social data.*

11 **Keywords:** R Shiny, Exploratory data analysis, Correlation network

---

Nr.	Code metadata description	Metadata
C1	Current code version	v3.5.4
C2	Permanent link to code/repository used for this code version	<a href="https://github.com/AntoineSoetewey/AssociationExplorer">https://github.com/AntoineSoetewey/AssociationExplorer</a>
C3	Permanent link to Reproducible Capsule	For example: <a href="https://codeocean.com/capsule/0270963/tree/v1xxx">https://codeocean.com/capsule/0270963/tree/v1xxx</a>
C4	Legal Code License	MIT License
C5	Code versioning system used	Git
C6	Software code languages, tools, and services used	R, R Shiny
C7	Compilation requirements, operating environments & dependencies	xxx
C8	If available link to developer documentation/manual	<a href="https://github.com/AntoineSoetewey/AssociationExplorer/tree/main/documentation">https://github.com/AntoineSoetewey/AssociationExplorer/tree/main/documentation</a> to do write doc xxx
C9	Support for questions or issues	<a href="https://github.com/AntoineSoetewey/AssociationExplorer/issues">https://github.com/AntoineSoetewey/AssociationExplorer/issues</a>

Table 1: Code metadata

## 12 Metadata

13 The metadata associated with the current version of the software is summa-  
14 rized in Table 1.

## 15 1. Motivation and significance

16 The growing availability of large, complex, and high-dimensional datasets in  
17 the social sciences and public policy domains offers unprecedented opportu-  
18 nities for insight but also presents significant challenges for exploration and  
19 interpretation, particularly for non-specialist audiences. Journalists, educa-  
20 tors, and engaged citizens often struggle to identify and interpret meaningful

---

\*Corresponding author.

*Email addresses:* [antoine.soetewey@uliege.be](mailto:antoine.soetewey@uliege.be) (Antoine Soetewey),  
[cedric.heuchenne@uclouvain.be](mailto:cedric.heuchenne@uclouvain.be) (Cédric Heuchenne), [arnaud.claes@uclouvain.be](mailto:arnaud.claes@uclouvain.be)  
(Arnaud Claes), [antonin.descampe@uclouvain.be](mailto:antonin.descampe@uclouvain.be) (Antonin Descampe)

relationships between variables without the aid of programming skills or formal statistical training. This barrier limits the broader societal impact of open data initiatives, which are designed to promote transparency, accountability, and informed public discourse.

To address this gap, we developed AssociationExplorer, a free, open-source R Shiny [1] application that enables intuitive and statistically grounded exploration of multivariate associations. The tool guides users through a visual journey of variable relationships by automatically computing appropriate bivariate association measures—Pearson’s  $r$ , Eta, and Cramer’s  $V$ —depending on variable types, and presenting the results in an interactive correlation network. Users can set thresholds for the strength of association and explore linked bivariate plots or tables with descriptive labels. This workflow supports transparent, reproducible, and non-technical exploratory data analysis (EDA).

Our software is particularly suited to survey-based datasets and public opinion studies. As an illustrative case, we apply AssociationExplorer to the European Social Survey (ESS), a cross-national survey that collects attitudinal, behavioral, and socio-demographic data across European countries. The tool allows users to uncover associations between trust in institutions, policy preferences, media usage, and demographic characteristics without any coding. This type of interactive analysis can empower journalists to build data-driven narratives, educators to teach statistical thinking, and citizens to explore evidence underlying public debate.

While several tools and libraries exist for correlation analysis (e.g., `corrr` [2] and `GGally` [4] in R [3], or Python packages like `seaborn` [6] and `pingouin` [5]), they typically require programming proficiency and focus on numerical associations. Other visualization tools such as Tableau or Power BI provide dashboards but often lack statistical rigor in bivariate association metrics or flexibility for categorical data. AssociationExplorer fills this gap by integrating statistical validity, accessibility, and interactivity in a single open-source web interface.

By lowering the technical barrier for statistical exploration, AssociationExplorer contributes to a more inclusive data culture and supports data-driven discovery in both academic and public-facing contexts.

## 2. Software description

### 2.1. Software architecture

The AssociationExplorer application is a web-based graphical user interface built with the R programming language using the Shiny framework. It adopts

59 a modular, reactive structure where data inputs and user selections dynami-  
 60 cally trigger updates to the visualizations and underlying computations. The  
 61 user interface is styled using the `bslib` package with a modern flat theme and  
 62 enhanced interactivity through `shinyjs` and `visNetwork`. The app is struc-  
 63 tured into distinct tabs: data upload, variable selection, correlation network  
 64 visualization, pairs plots, and a help section.  
 65 Upon upload, the dataset is preprocessed to exclude variables with zero vari-  
 66 ance, as these variables do not vary across observations and therefore cannot  
 67 contribute to meaningful associations or visualizations. Removing them helps  
 68 reduce noise and ensures that only informative variables are included in the  
 69 analysis. Optionally, the user can provide a variable description file, which  
 70 is integrated and used to annotate visual elements. The backend computes  
 71 association measures tailored to the variable types: Pearson’s  $r$  for numeric  
 72 pairs, Cramer’s  $V$  for categorical pairs, and the correlation ratio ( $\eta$ ) for  
 73 mixed pairs. Associations are filtered using user-defined thresholds and rep-  
 74 resented in a correlation network and complementary bivariate plots. The  
 75 app handles both CSV and Excel files and supports large datasets of up to  
 76 100 MB.

## 77 *2.2. Software functionalities*

78 The major functionalities of the AssociationExplorer application include:

- 79 • **Data ingestion and cleaning:** The app supports CSV and Excel  
 80 files. It automatically removes variables with only one unique value,  
 81 as they lack variability and cannot contribute to association analyses.  
 82 Additionally, it can optionally integrate user-supplied descriptions of  
 83 variables, which are used to enhance the clarity and interpretability of  
 84 visualizations, particularly for non-technical users.
- 85 • **Variable selection interface:** Users can interactively choose which  
 86 variables to explore. When a description file is provided, a summary  
 87 table links variable names to their descriptions.
- 88 • **Dynamic association filtering:** The app computes pairwise associ-  
 89 ation measures between all selected variables, using a method tailored  
 90 to the types of variables involved:
  - 91 – For pairs of numeric variables, the app calculates Pearson’s corre-  
 92 lation coefficient ( $r$ ), and retains the association if the coefficient  
 93 of determination ( $R^2$ ) exceeds a user-defined threshold:

$$R^2 = r^2$$

94 where  $r = \text{cor}(X, Y)$ .

- 95 – For pairs of categorical variables, it computes Cramer’s V, a nor-  
 96 malized measure of association derived from the chi-squared statis-  
 97 tic:

$$V = \sqrt{\frac{\chi^2}{n \cdot \min(k - 1, r - 1)}} \quad (1)$$

98 where  $\chi^2$  is the chi-squared statistic,  $n$  is the total number of ob-  
 99 servations, and  $k, r$  are the number of categories in each variable.

- 100 – For mixed pairs (one numeric and one categorical variable), the  
 101 app computes the correlation ratio ( $\eta$ ), which quantifies how much  
 102 of the variance in the numeric variable is explained by the grouping  
 103 structure of the categorical variable. It is defined as:

$$\eta = \sqrt{\frac{\text{SS}_{\text{between}}}{\text{SS}_{\text{total}}}} \quad (2)$$

104 where:

- 105 -  $\text{SS}_{\text{total}}$  is the *total sum of squares* of the numeric variable:

$$\text{SS}_{\text{total}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

106 with  $y_i$  the observed numeric values and  $\bar{y}$  their overall mean.

- 107 -  $\text{SS}_{\text{between}}$  is the *between-group sum of squares*, computed as:

$$\text{SS}_{\text{between}} = \sum_{g=1}^G n_g (\bar{y}_g - \bar{y})^2$$

108 where  $G$  is the number of groups (categories),  $n_g$  is the number  
 109 of observations in group  $g$ ,  $\bar{y}_g$  is the group mean, and  $\bar{y}$  is the  
 110 overall mean.

111 This formulation captures the proportion of the total variance in  
 112 the numeric variable that can be attributed to differences between  
 113 the categorical groups. A pair is retained only if  $\eta^2$  exceeds the  
 114 numeric threshold defined by the user.

115 Each association is retained only if its corresponding strength metric–  
 116  $R^2$ ,  $\eta^2$ , or Cramer’s V–exceeds the threshold set by the user. These

117 thresholds can be adjusted interactively through the interface, and  
118 the filtering process is reactive: updates to the thresholds immediately  
119 propagate to the network and bivariate visualizations. This allows users  
120 to dynamically control the sensitivity of the association analysis and  
121 focus on relationships of substantive interest.

122 • **Interactive correlation network:** The filtered associations are dis-  
123 played as an interactive graph where nodes represent variables and  
124 edges represent associations. Edge thickness and length reflect the  
125 strength of the association: stronger associations are shown with thicker  
126 and shorter edges, whereas weaker associations are displayed with thin-  
127 ner and longer edges. This dual visual representation helps users quickly  
128 identify the most meaningful relationships in the network. Variable  
129 descriptions are displayed when the user hovers over a node in the net-  
130 work, allowing for quick access to additional context without cluttering  
131 the visualization.

132 xxx add more details, check visnetwork documentation for example.

133 • **Bivariate visualization of variable pairs:** For each variable pair  
134 exceeding the threshold:

- 135 – Scatter plots with linear regression lines are shown for numeric  
136 pairs, helping visualize the direction and strength of the relation-  
137 ship.
- 138 – Colored contingency tables with marginal sums are shown for cat-  
139 egorical pairs, where cell background colors vary in intensity ac-  
140 cording to the frequency of observations, using a blue gradient to  
141 highlight higher counts.
- 142 – Mean plots are shown for numeric-categorical pairs, with bars or-  
143 dered by mean value to make it easy to compare and rank cate-  
144 gories based on the quantitative variable.

145 Confidence intervals for the regression lines and standard errors in the  
146 mean plots are intentionally omitted to maintain a clean, uncluttered  
147 visualization that prioritizes ease of interpretation. Mean plots were  
148 selected over boxplots to avoid overwhelming non-expert users with  
149 distributional information, focusing instead on clear, accessible insights  
150 about average group differences.

151 • **Accessibility and user guidance:** A dedicated help section explains  
152 each step, allowing users with a limited statistical background to inter-  
153 actively explore their data.

154 *2.3. Sample code snippets analysis*

155 Below is a representative snippet from the application showing how the soft-  
156 ware selects the appropriate association measure depending on the types of  
157 the variable pair and filters associations based on user thresholds:

```
158 # Numeric vs numeric case
159 if (is_num1 && is_num2) {
160     ...
161     r <- cor(x, y, use = "complete.obs")
162     cor_val <- ifelse(r^2 >= threshold_num, r, 0)
163     cor_type <- "Pearson's r"
164
165 # Categorical vs categorical case
166 } else if (!is_num1 && !is_num2) {
167     ...
168     tbl <- table(x, y)
169     ...
170     n_obs <- sum(tbl)
171     df_min <- min(nrow(tbl) - 1, ncol(tbl) - 1)
172     if (df_min > 0) {
173         v_cramer <- sqrt(chi$statistic / (n_obs * df_min))
174         cor_val <- ifelse(v_cramer >= threshold_cat,
175                           v_cramer, 0)
176         cor_type <- "Cramer's V"
177     }
178
179 # Mixed case (numeric vs categorical)
180 } else {
181     ...
182     means_by_group <- tapply(num_var, cat_var,
183                               mean, na.rm = TRUE)
184     overall_mean <- mean(num_var, na.rm = TRUE)
185     n_groups <- tapply(num_var, cat_var, length)
186     bss <- sum(n_groups * (means_by_group - overall_mean)^2,
187               na.rm = TRUE)
188     tss <- sum((num_var - overall_mean)^2, na.rm = TRUE)
189
190     if (tss > 0) {
191         eta <- sqrt(bss / tss)
192         cor_val <- ifelse(eta^2 >= threshold_num, eta, 0)
193         cor_type <- "Eta"
```

194 }  
195 }

196 This conditional structure ensures that the correct statistical method is ap-  
197 plied for each type of variable pair, supporting a robust and interpretable  
198 exploration of associations.  
199 xxx add screenshots

### 200 3. Illustrative examples

201 *Provide at least one illustrative example to demonstrate the major functions*  
202 *of your software/code.*

203 **Optional:** *you may include one explanatory video or screencast that will ap-*  
204 *pear next to your article, in the right hand side panel. Please upload any video*  
205 *as a single supplementary file with your article. Only one MP4 formatted,*  
206 *with 150MB maximum size, video is possible per article. Recommended video*  
207 *dimensions are 640 x 480 at a maximum of 30 frames / second. Prior to sub-*  
208 *mission please test and validate your .mp4 file at [http://elsevier-apps.](http://elsevier-apps.sciiverse.com/GadgetVideoPodcastPlayerWeb/verification)*  
209 *sciiverse.com/GadgetVideoPodcastPlayerWeb/verification* . *This tool*  
210 *will display your video exactly in the same way as it will appear on ScienceDi-*  
211 *rect.*

### 212 4. Impact

213 *This is the main section of the article and reviewers will weight it appropri-*  
214 *ately. Please indicate:*

- 215 • *Any new research questions that can be pursued as a result of your*  
216 *software.*
- 217 • *In what way, and to what extent, your software improves the pursuit of*  
218 *existing research questions.*
- 219 • *Any ways in which your software has changed the daily practice of its*  
220 *users.*
- 221 • *How widespread the use of the software is within and outside the in-*  
222 *tended user group (downloads, number of users if your software is a*  
223 *service, citable publications, etc.).*
- 224 • *How the software is being used in commercial settings and/or how it*  
225 *has led to the creation of spin-off companies.*

226 *Please note that points 1 and 2 are best demonstrated by references to citable*  
227 *publications.*



## 228 5. Conclusions

## 229 References

- 230 [1] Chang, W., Cheng, J., Allaire, J., Sievert, C., Schloerke, B., Xie, Y.,  
231 Allen, J., McPherson, J., Dipert, A., and Borges, B. (2024). *shiny: Web*  
232 *Application Framework for R*. R package version 1.9.1.
- 233 [2] Kuhn, M., Jackson, S., and Cimentada, J. (2022). *corrr: Correlations in*  
234 *R*. R package version 0.4.4.
- 235 [3] R Core Team (2024). *R: A Language and Environment for Statistical*  
236 *Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- 237 [4] Schloerke, B., Cook, D., Larmarange, J., Briatte, F., Marbach, M.,  
238 Thoen, E., Elberg, A., and Crowley, J. (2024). *GGally: Extension to*  
239 *'ggplot2'*. R package version 2.2.1.
- 240 [5] Vallat, R. (2018). Pingouin: statistics in python. *Journal of Open Source*  
241 *Software*, 3(31):1026.
- 242 [6] Waskom, M. L. (2021). seaborn: statistical data visualization. *Journal*  
243 *of Open Source Software*, 6(60):3021.