

The Strucplot Framework—Visualizing Multi-way Contingency Tables

by David Meyer, Achim Zeileis, and Kurt Hornik

1 Framework Overview

The strucplot framework in the R package **vcd**, used for visualizing multi-way contingency tables, integrates techniques such as mosaic displays, association plots, and sieve plots. The main idea is to visualize the tables’ cells arranged in rectangular form. For multi-way tables, the variables are nested into rows and columns using recursive conditional splits, given the margins. The result is a ‘flat’ representation that can be visualized in ways similar to a two-dimensional table. This principle defines a class of conditional displays which allows for granular control of graphical appearance aspects, including:

- the content of the tiles
- the split direction for each dimension
- the graphical parameters of the tiles’ content
- the spacing between the tiles
- the labeling of the tiles

This document gives an introduction to the framework, whereas labeling and shading issues are described in separate vignettes.

The strucplot framework is highly modularized: Figure 1 shows the hierarchical relationship between the various components. On the lowest level, there are several groups of workhorse and parameter functions that directly or indirectly influence the final appearance of the plot. These are examples of ‘graphical appearance control’ (‘grapcon’) functions. They are created by generating functions (‘grapcon generators’), allowing flexible parameterization and extensibility (Figure 1 only shows the generators). The first part of the generator names (*group_foo()*) reflects the group they belong to (strucplot core, labeling, legend, shading, or spacing). The workhorse functions (created by *struc_foo()*, *labeling_foo()*, and *legend_foo()*) directly produce graphical output (“add ink to the canvas”), whereas the parameter functions (created by *spacing_foo()* and *shading_foo()*) compute graphical parameters used by the others. The grapcon functions returned by *struc_foo()* implement the core functionality, creating the tiles and their content. On the second level of the framework, a suitable combination of the low-level grapcon functions (or, alternatively, corresponding generating functions) is passed as “hyperparameters” to *strucplot()*. This central function sets up the graphical layout using grid viewports (see Figure 2), and coordinates the specified core, labeling, shading, and spacing functions to produce the plot. On the third level, we provide several convenience functions such as *mosaic()*, *sieve()*, *assoc()*, and *doubledecker()* which interface *strucplot()* through sensible parameter defaults and support for model formulas. Finally, on the fourth level, there are ‘related’ **vcd** functions (such as *cotabplot()* and the *pairs()* methods for table objects) arranging collections of plots of the strucplot framework into more complex displays (e.g., by means of panel functions).

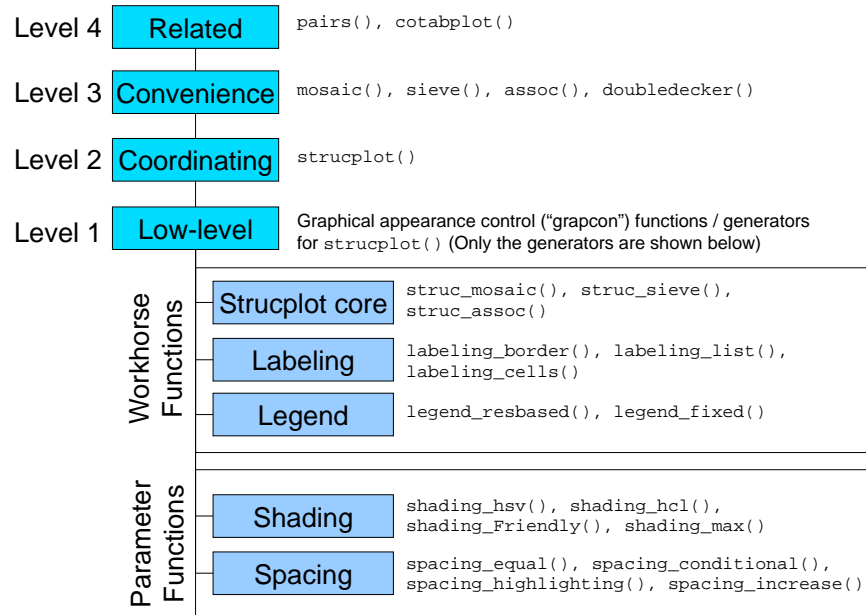


Figure 1: Components of the strucplot framework.

2 Mosaic, Association, and Sieve Plots

As an example, consider the ‘HairEyeColor’ data containing two polytomous variables (hair and eye color), as well as one (artificial) dichotomous variable (sex, i.e., gender). The ‘flattened’ contingency table can be obtained using the `structable()` function (quite similar to `fTable()` in base R, but allowing the specification of split directions):

```
> (hec <- structable(Eye ~ Sex + Hair, data = HairEyeColor))
```

		Eye			
		Brown	Blue	Hazel	Green
Sex	Hair				
	Male				
	Black	32	11	10	3
	Brown	38	50	25	15
Female	Red	10	10	7	7
	Blond	3	30	5	8
	Black	36	9	5	2
	Brown	81	34	29	14
	Red	16	7	7	7
	Blond	4	64	5	8

Let us first visualize the contingency table by means of a mosaic plot (Hartigan and Kleiner, 1984) which is basically an area-proportional visualization of (typically, observed) frequencies, composed of tiles (corresponding to the cells) created by recursive vertical and horizontal splits of a square. Thus, the area of each tile is proportional to the corresponding cell entry *given* the dimensions of previous splits. Figure 3 depicts the effect of

```
> mosaic(hec)
```

equivalent to

main			
[A]	margin_top	[B]	[E]
margin_left	plot	margin_right	legend
[C]	margin_bottom	[D]	[F]
sub			

Figure 2: Viewport layout for strucplot displays with their names. [A] = “corner_top_left”, [B] = “corner_top_right”, [C] = “corner_bottom_left”, [D] = “corner_bottom_right”, [E] = “legend_top”, [F] = “legend_sub”.

```
> mosaic(~Sex + Eye + Hair, data = HairEyeColor)
```

The small bullets indicate zero entries in the corresponding cell. Note that in contrast to, e.g., `mosaicplot()` in base R, the default split direction and level ordering in all strucplot displays correspond to the textual representation. It is also possible to visualize the expected values instead of the observed values (see Figure 4):

```
> mosaic(hec, type = "expected")
```

In order to compare observed and expected values, a sieve plot (Riedwyl and Schüpbach, 1994) could be used (see Figure 5):

```
> sieve(hec)
```

Alternatively, we can directly inspect the residuals. The Pearson residuals (standardized deviations of observed from expected values) are preferably visualized using association plots (Cohen, 1980). In contrast to `assocplot()` in base R, `vcd`’s `assoc()` function scales to more than two variables (see Figure 6):

```
> assoc(hec, compress = FALSE)
```

The `compress` argument keeps distances between tiles equal for better comparison.

For both mosaic plots and association plots, the splitting of the tiles can be controlled using the `split_vertical` argument (default: alternating splits starting with a vertical one).

```
> mosaic(hec, split_vertical = c(TRUE, FALSE, TRUE),
+       labeling_args = list(abbreviate = c(Eye = 3)))
```

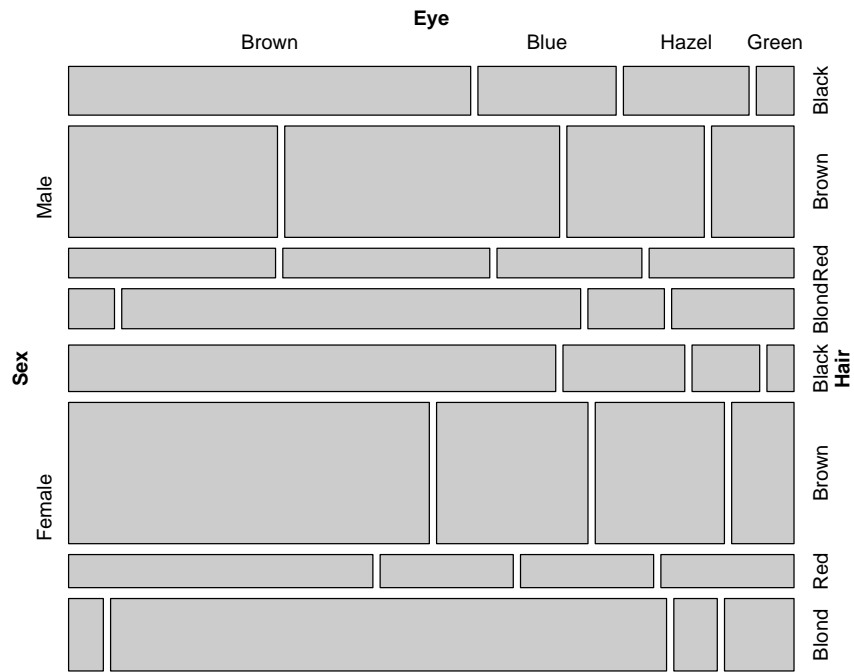


Figure 3: Mosaic plot for the 'HairEyeColor' data.

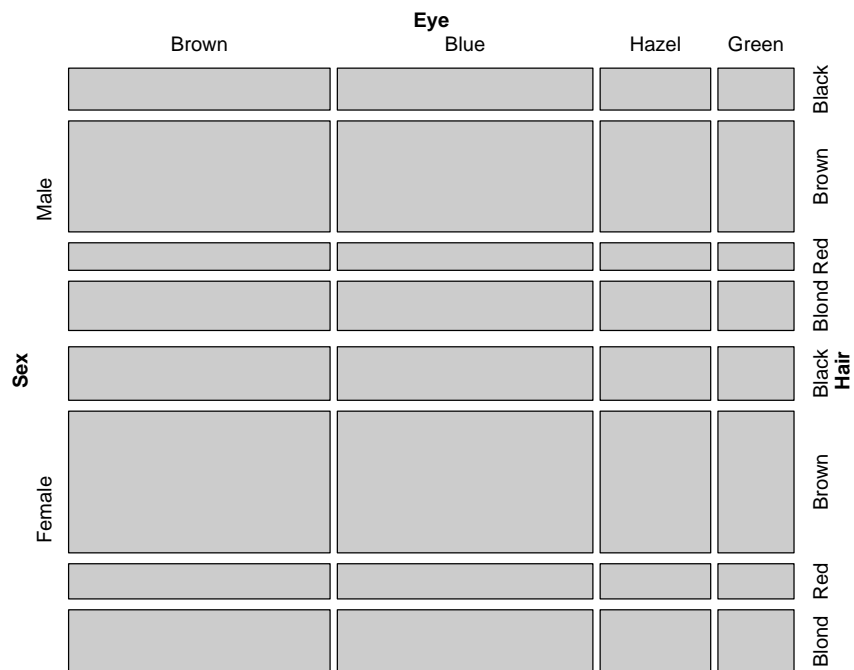


Figure 4: Mosaic plot for the 'HairEyeColor' data (expected values).

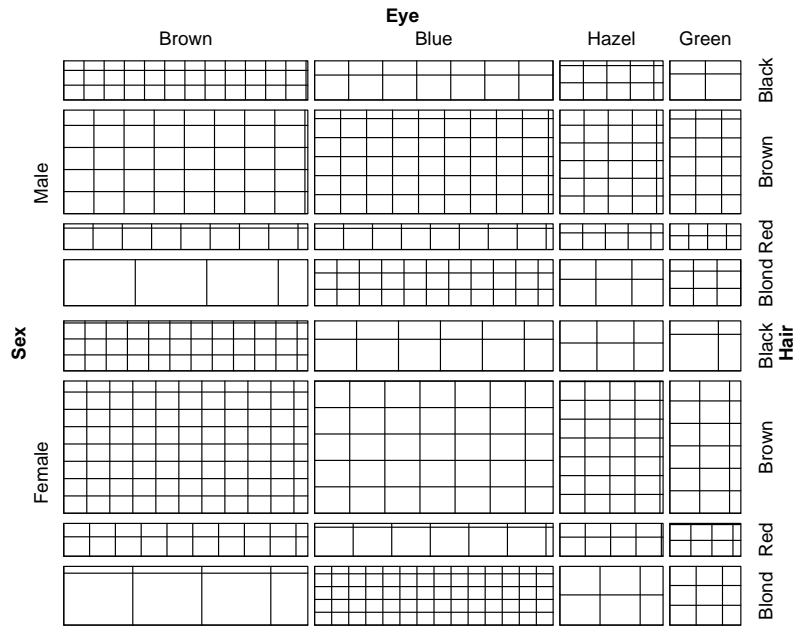


Figure 5: Sieve plot for the ‘HairEyeColor’ data visualizing simultaneously observed and expected values.

For compatibility with `mosaicplot()` in base R, the `mosaic()` function also allows the use of a “direction” argument taking a vector of “h” and “v” characters (see Figure 7):

```
> mosaic(hec, direction = c("v", "h", "v"))
```

By a suitable combination of splitting, spacing, and labeling settings, the functions provided by the `strucplot` framework can be customized in a quite flexible way. For example, `doubledecker()` is simply a wrapper for `mosaic()`, setting the right defaults. Figure 8 shows a `doubledecker` plot of the ‘Titanic’ data, explaining the probability of survival (‘survived’) by age, given sex, given class. It is created by:

```
> doubledecker(Titanic)
```

equivalent to:

```
> doubledecker(Survived ~ Class + Sex + Age, data = Titanic)
```

3 Conditional and partial views

So far, we have visualized full tables. Subtables can be selected in a similar way than for objects of class `table` using indexing. Note, however, that subsetting of `structable` objects is more restrictive because of their inherent conditional structure. Since the variables on both the row and the columns side are nested, subsetting is only possible “outside-in”, that is, indexing operates on blocks defined by the variable levels:

```
> hec
```

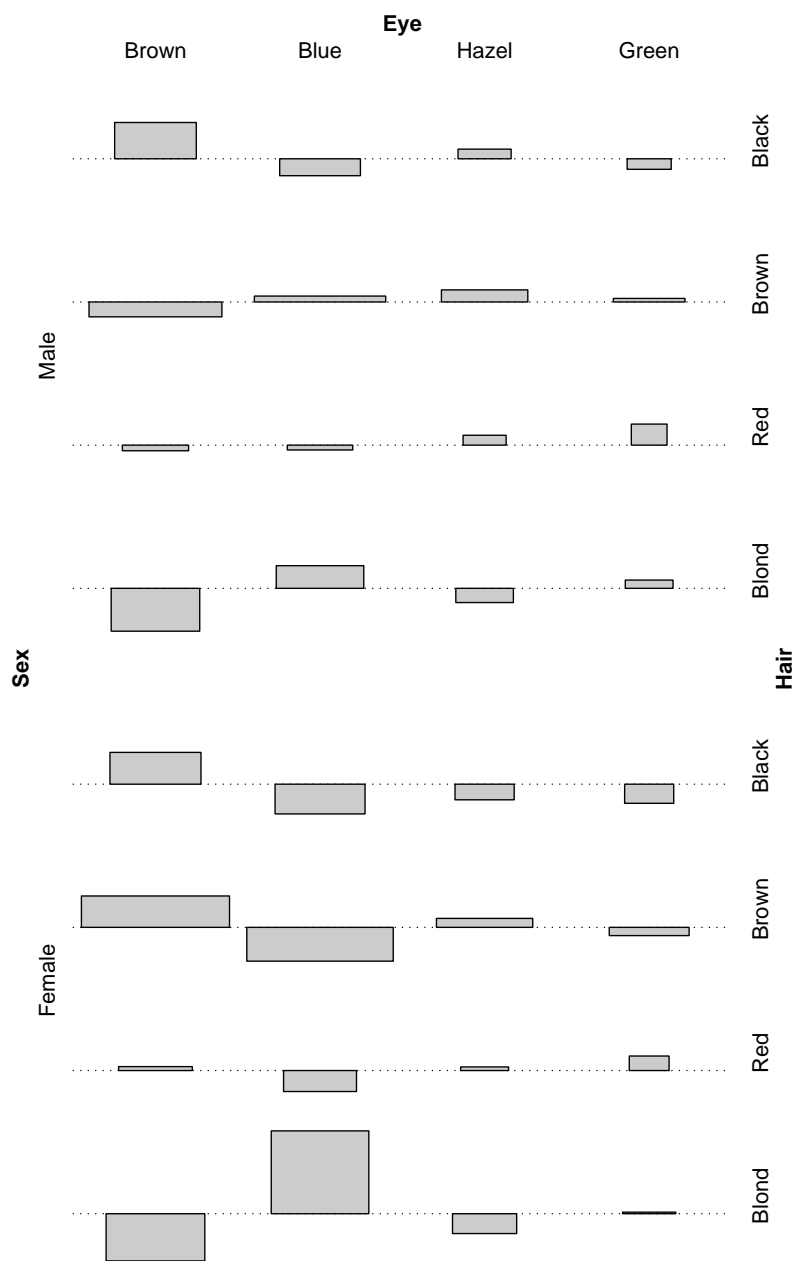


Figure 6: Association plot for the 'HairEyeColor' data.

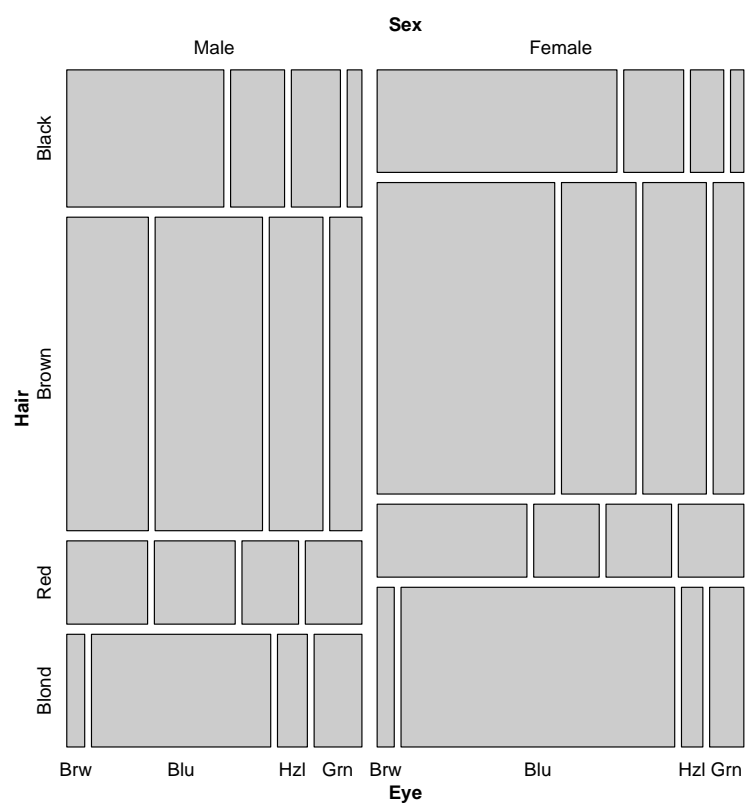


Figure 7: Mosaic plot for the ‘HairEyeColor’ data—alternative splitting.

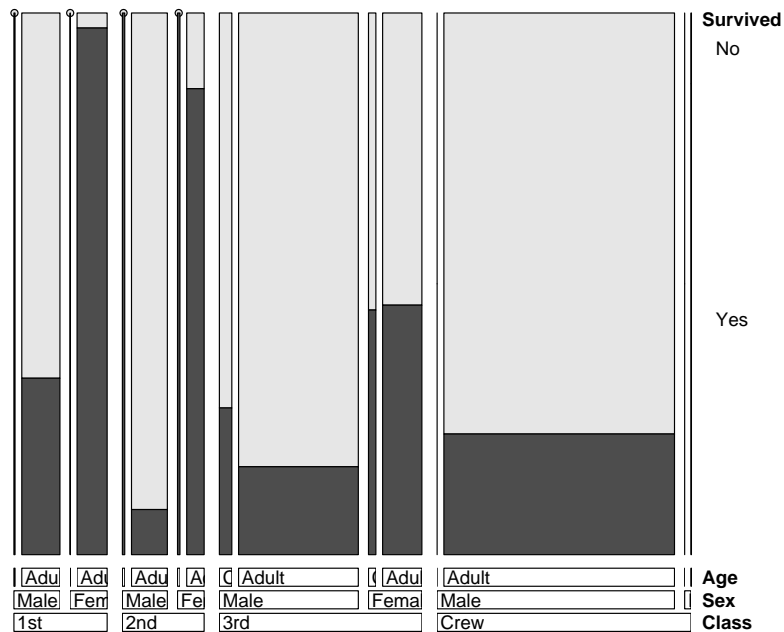


Figure 8: Doubled-decker plot for the ‘Titanic’ data.

		Eye Brown Blue Hazel Green			
Sex	Hair				
Male	Black	32	11	10	3
	Brown	38	50	25	15
	Red	10	10	7	7
	Blond	3	30	5	8
Female	Black	36	9	5	2
	Brown	81	34	29	14
	Red	16	7	7	7
	Blond	4	64	5	8

```
> hec["Male", ]
```

		Eye Brown Blue Hazel Green			
Sex	Hair				
Male	Black	32	11	10	3
	Brown	38	50	25	15
	Red	10	10	7	7
	Blond	3	30	5	8

```
> hec["Male", c("Brown", "Green")]
```

		Eye Brown Green	
Sex	Hair		
Male	Black	32	3
	Brown	38	15
	Red	10	7
	Blond	3	8

Conditioning on levels (i.e., choosing a table subset for fixed levels of the conditioning variable(s)) is done using the `[[` operator. Here again, the sequence of conditioning levels is restricted by the hierarchical structure of the `structable` object:

```
> hec[["Male", ]]

      Eye Brown Blue Hazel Green
Hair
Black      32   11    10     3
Brown      38   50    25    15
Red        10   10     7     7
Blond       3   30     5     8
```

```
> hec[[c("Male", "Red"), ]]

      Eye Brown Blue Hazel Green
      10   10     7     7
```

```
> hec[["Male", "Green"]]

Hair
Black   3
Brown  15
Red     7
Blond   8
```

Now, there are several ways for visualizing conditional independence structures. The “brute force” method is to draw separate plots for the strata. The following example compares the association between hair and eye color, given gender, by using subsetting on the flat table and **grid**’s viewport framework to visualize the two groups besides each other:

```
> pushViewport(viewport(layout = grid.layout(ncol = 2)))
> pushViewport(viewport(layout.pos.col = 1))
> mosaic(hec[["Male"]], margins = c(left = 2.5, top = 2.5,
+   0), sub = "Male", newpage = FALSE)
> popViewport()
> pushViewport(viewport(layout.pos.col = 2))
> mosaic(hec[["Female"]], margins = c(top = 2.5, 0), sub = "Female",
+   newpage = FALSE)
> popViewport(2)
```

Note the use of the `margins` argument: it takes a vector with up to four values whose unnamed components are recycled, but “overruled” by the named arguments. Thus, in the example, only the top margin is set to 2 lines, and all other to 0. This idea applies to almost all vectorized arguments in the `strucplot` framework (with `split_vertical` as a prominent exception).

Since `mosaic` displays are “conditional plots” by definition, we can also use one single `mosaic` for stratified plots. The formula interface of `mosaic()` allows the specification of conditioning variables (see Figure 10):

```
> mosaic(~Hair + Eye | Sex, data = hec, split_vertical = TRUE,
+   keep_aspect_ratio = FALSE)
```

The effect of using this kind of formula is that conditioning variables are permuted ahead of the the conditioned variables in the table, and that `spacing_conditional()` is used as default to better distinguish conditioning from conditioned dimensions. This spacing uses equal space between tiles

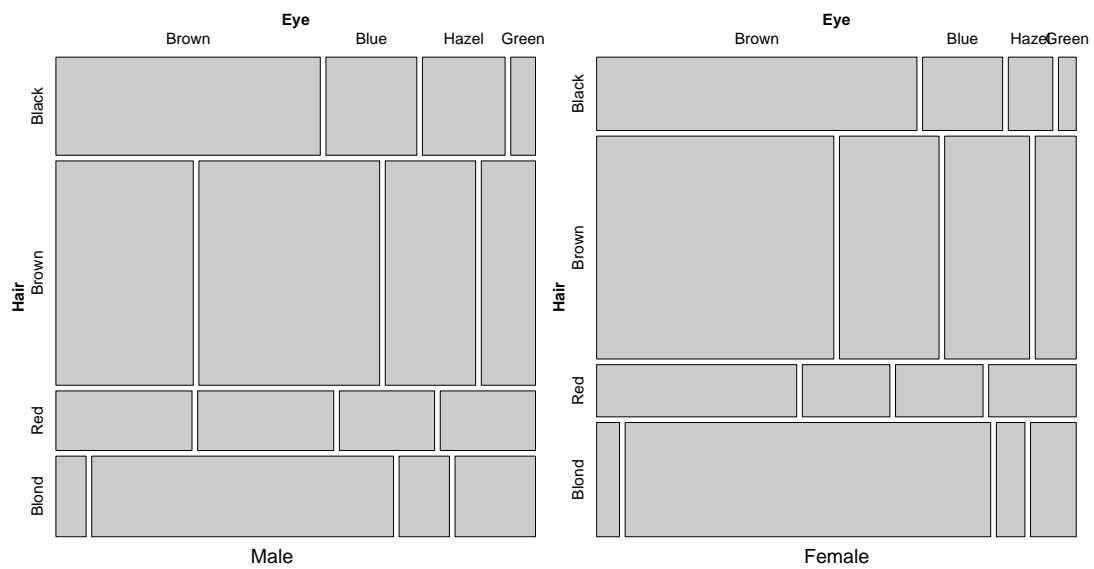


Figure 9: Distribution of hair and eye color, given gender.



Figure 10: Mosaic plot for conditional independence structures.

of conditioned variables, and increasing space between tiles of conditioning variables. In addition, we release the fixed aspect ratio to get less distorted margins.

The `cotabplot()` function does a much better job on this task: it arranges stratified strucplot displays in a lattice-like layout, conditioning on variable *levels*. The plot in Figure 11 shows hair and eye color, given sex:

```
> cotabplot(~Hair + Eye | Sex, data = hec, panel_args = list(margins = 3),
+           labeling = labeling_left(clip = FALSE))
```

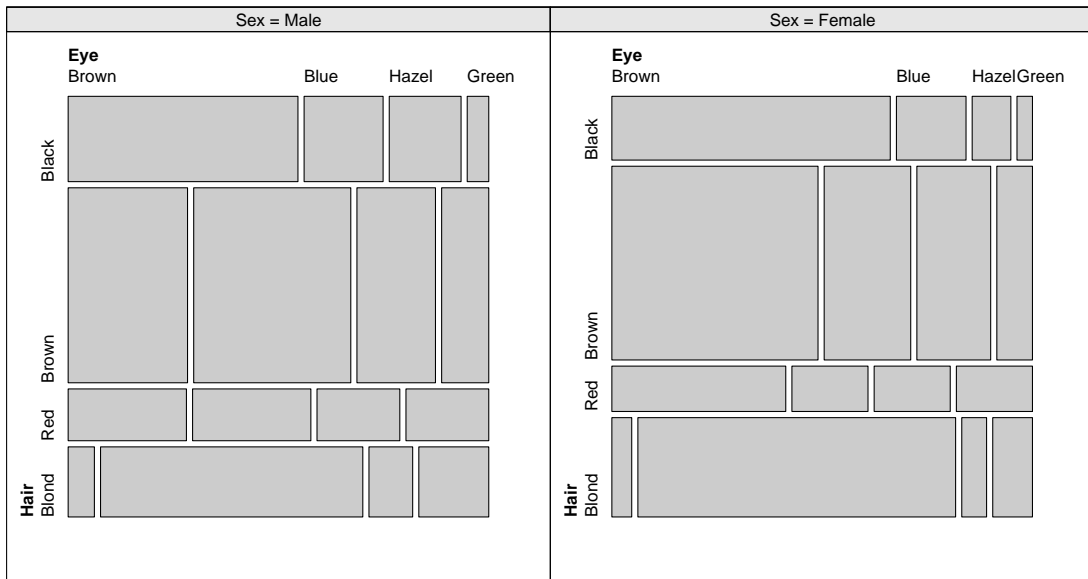


Figure 11: Conditional table plot for the ‘HairEyeColor’ data.

The `labeling_args` argument modifies the labels’ appearance: here, to be left-aligned and unclipped (see the separate vignette: “Labeling in the Strucplot Framework” for detailed information).

Another high-level function for visualizing conditional independence models are the `pairs()` methods for `table` and `strutable` objects. In contrast to `cotabplot()` which conditions on variables, the `pairs()` methods create pairwise views of the table. The function produces, by default, a plot matrix having strucplot displays in the off-diagonal panels, and the variable names (optionally, with univariate statistics) in the diagonal cells. Figure 12 shows a `pairs` display with mosaic plots visualizing mutual independence in the lower triangle, association plots for the same in the upper triangle, and bar charts in the diagonal.

```
> pairs(hec, lower_panel = pairs_assoc, space = 0.3, diag_panel_args = list(rot = -45,
+   just_leveltext = c("left", "bottom")))
```

(The labels of the variables are to be read from left to right and from top to bottom.) In plots produced by `pairs()`, each panel’s row and column define two variables X and Y used for the specification of four different types of independence: pairwise, total, conditional, and joint. The pairwise mosaic matrix shows bivariate marginal relations between X and Y , collapsed over all other variables. The total independence mosaic matrix shows mosaic plots for mutual independence, i.e., for marginal and conditional independence among all pairs of variables. The conditional independence mosaic matrix shows mosaic plots for marginal independence of X and Y , given all other variables. The joint independence mosaic matrix shows mosaic plots for joint independence of all pairs (X, Y) of variables from the others.

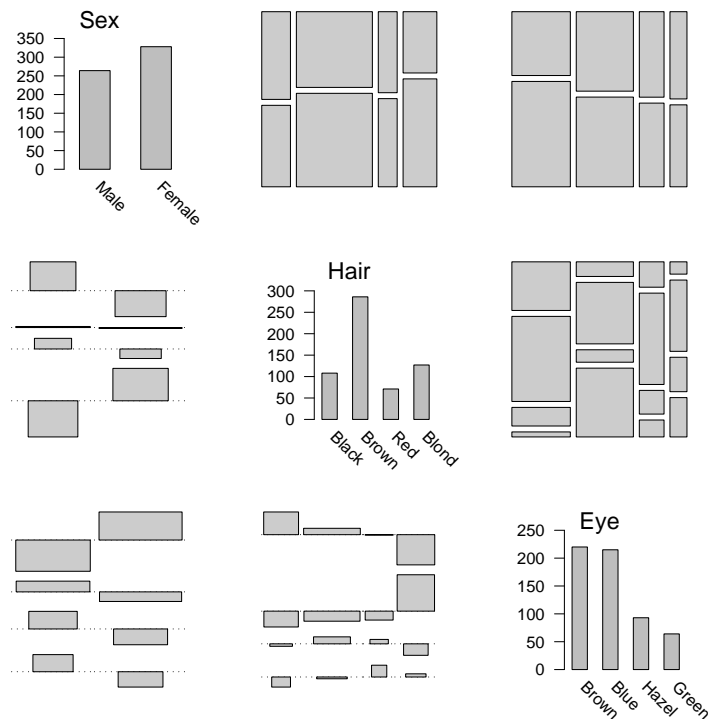


Figure 12: Pairs plot for the ‘HairEyeColor’ data.

Since the matrix is symmetric, the upper and lower parts can independently be used to display different types of independence models, or different strucplots displays (mosaic, association, or sieve plots). The available panel functions (`pairs_assoc()`, `pairs_mosaic()`, and `pairs_sieve()`) are simple wrappers to `assoc()`, `mosaic()`, and `sieve()`, respectively. Obviously, seeing patterns in strucplot matrices becomes increasingly difficult with higher dimensionality. Therefore, this plot is typically used with a suitable residual-based shading (described in the vignette on “Colors and Residual-based Shadings in the Strucplot Framework”).

4 Interactive plot modifications

All strucplot core functions are supposed to produce conditional hierarchical plots by the means of nested viewports, corresponding to the provided splitting information. Thus, at the end of the plotting, each tile is associated with a particular viewport. Each of those viewports has to be conventionally named, enabling other strucplot modules, in particular the labeling functions, to access specific tiles after they have been plotted. The naming convention for the viewports is:

`cell:Variable1=Level1,Variable2=Level2 ...`

Clearly, these names depend on the splitting. The following example shows how to access parts of the plot after it has been drawn (see Figure 13):

```
> mosaic(~Hair + Eye, data = hec, pop = FALSE)
> seekViewport("cell:Hair=Blond")
> grid.rect(gp = gpar(col = "red", lwd = 4))
```

```
> seekViewport("cell:Hair=Blond, Eye=Blue")
> grid.circle(r = 0.2, gp = gpar(fill = "cyan"))
```

Note that the viewport tree is removed by default. Therefore, the `pop` argument has to be set to `FALSE` when viewports shall be accessed.

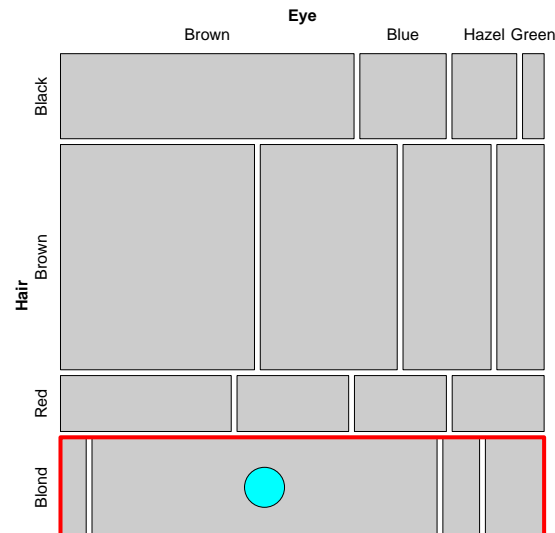


Figure 13: Adding elements to a mosaic plot after drawing.

In addition to the viewports, the main graphical elements get names following a similar construction method. This allows to change graphical parameters of plot elements *after* the plotting (see Figure 14):

```
> assoc(Eye ~ Hair, data = hec, pop = FALSE)
> getNames()[1:6]

[1] "GRID.lines.4817"      "rect:Hair=Black, Eye=Brown"
[3] "GRID.lines.4818"      "rect:Hair=Black, Eye=Blue"
[5] "GRID.lines.4819"      "rect:Hair=Black, Eye=Hazel"

> grid.edit("rect:Hair=Blond, Eye=Blue", gp = gpar(fill = "red"))
```

References

- Cohen A (1980). “On the Graphical Display of the Significant Components in a Two-Way Contingency Table.” *Communications in Statistics—Theory and Methods*, **A9**, 1025–1041.
- Hartigan J, Kleiner B (1984). “A mosaic of television ratings.” *The American Statistician*, **38**, 32–35.
- Riedwyl H, Schüpbach M (1994). “Parquet diagram to plot contingency tables.” In F Faulbaum (ed.), “Softstat ’93: Advances in Statistical Software,” pp. 293–299. Gustav Fischer, New York.

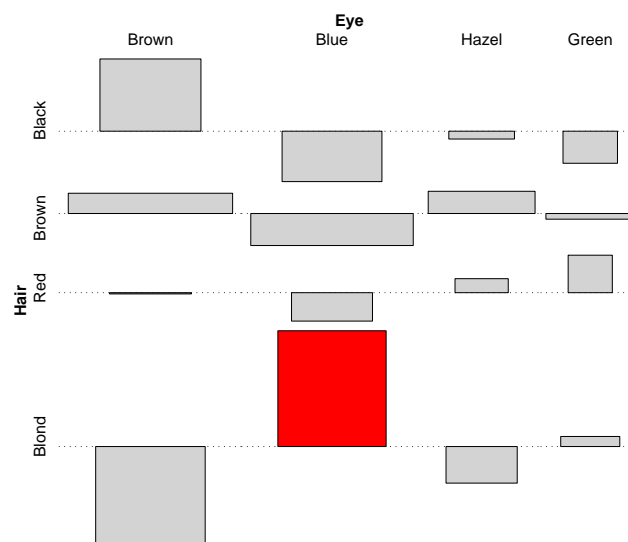


Figure 14: Changing graphical parameters of elements after drawing.