



Data Inspection Using Biplots

Ulrich Kohler

kohler@wz-berlin.de

Wissenschaftszentrum Berlin

Plan of the presentation

- ⑥ History
- ⑥ Interpretation
- ⑥ The math
- ⑥ Biplot-Types
- ⑥ Two more options

History

`biplot` has been available on SSC since Stata 5. After arrival of Stata 8 I have revisited `biplot` and made several changes (old version still works under version control).

- ⑥ Use of the new graph engine
- ⑥ Allowing for weights for JK-Biplots
- ⑥ New option `rv` for “compositional data”
- ⑥ New option `mahalanobis`
- ⑥ New option `subpop()`
- ⑥ Change of some default settings

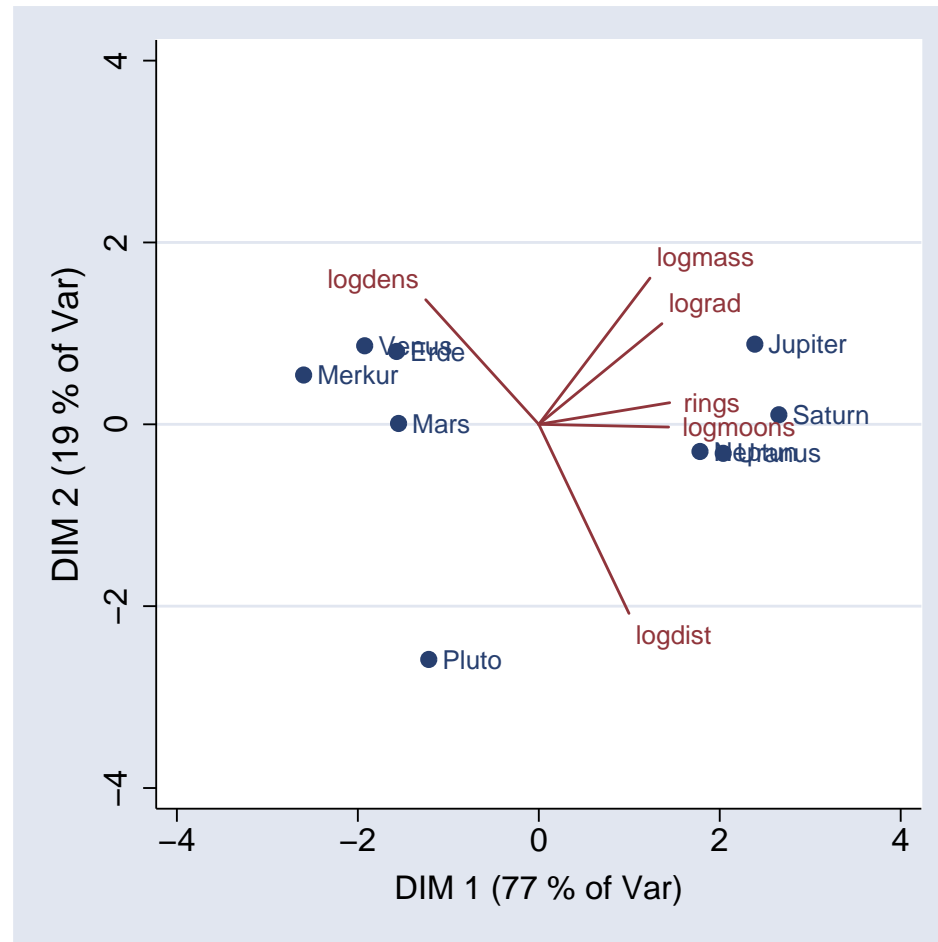
Interpretation

Biplots show the following quantities of a data matrix in one display:

- ⑥ standard deviations of variables
- ⑥ correlations between variables
- ⑥ values of observations on variables
- ⑥ distances between observations in the multidimensional space

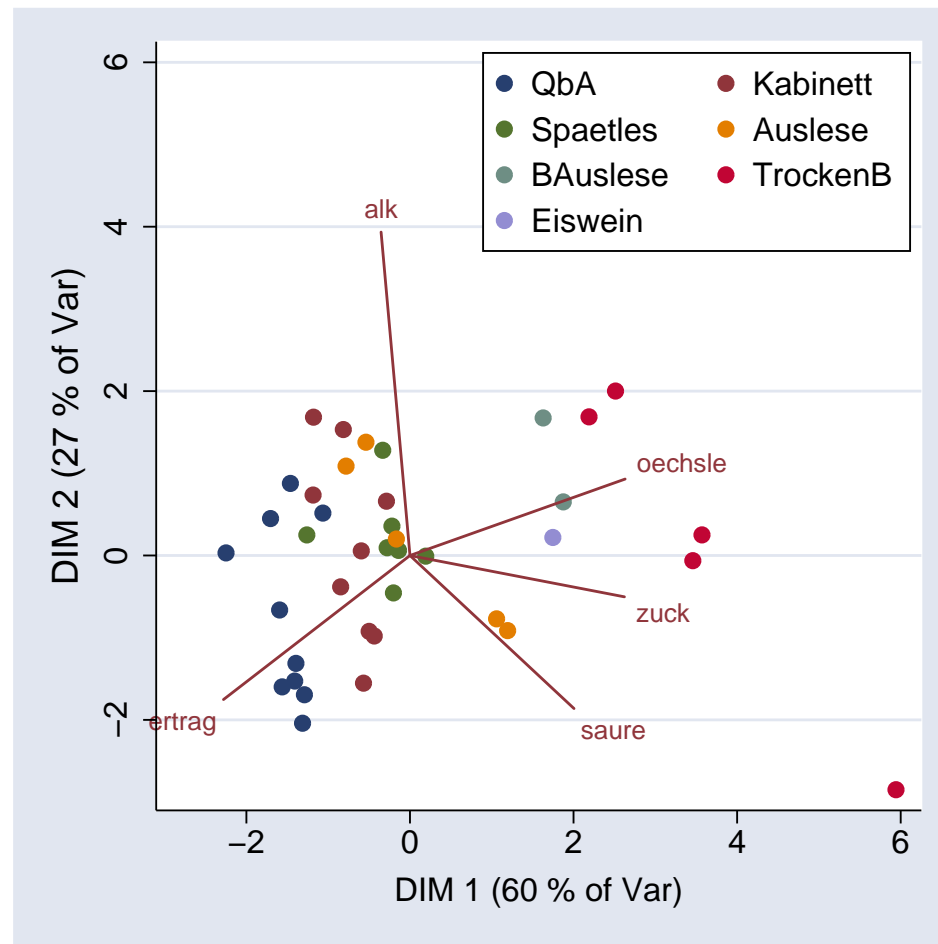
Interpretation

- `biplot ring-logmoons, mlabel(planet)`



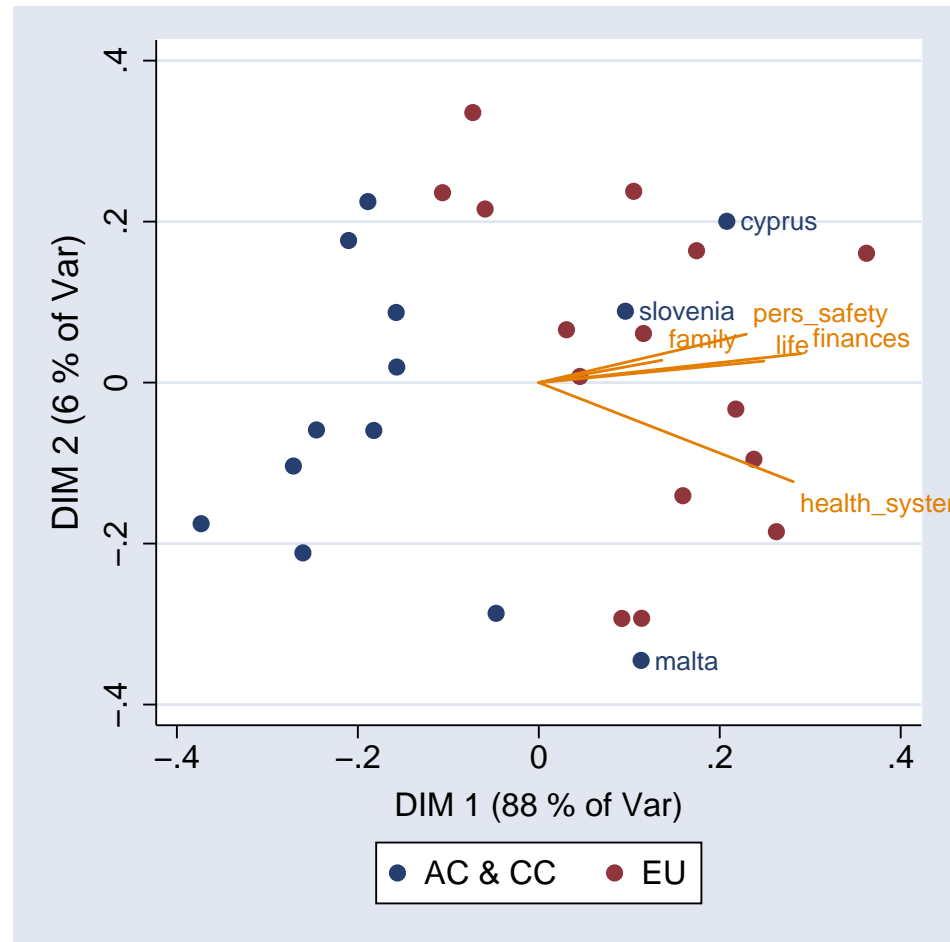
Interpretation

```
. biplot ..., subpop(praed) legend(ring(0)  
pos(1))
```



Interpretation

```
biplot ... , gh cov subpop(eu, mlab(label))
```



The Math

Let \mathbf{Y} be a $n \times k$ matrix holding the data. One can decompose \mathbf{Y} with a *singular value decomposition* (SVD) into

$$\underset{n \times k}{\mathbf{Y}} = \underset{n \times k}{\mathbf{U}} \underset{k \times k}{\mathbf{L}} \underset{k \times k}{\mathbf{V}'} \quad (1)$$

where \mathbf{L} contains the *Eigenvalues*.

From the SVD results the 2×2 matrix $\underline{\mathbf{L}}$ is formed, which contains the two elements of \mathbf{L} with the highest Eigenvalues. The $n \times 2$ matrix $\underline{\mathbf{U}}$ and the $k \times 2$ matrix $\underline{\mathbf{V}}$ are formed by choosing those columns from \mathbf{V} and \mathbf{U} which correspond to the highest Eigenvalues.

The Math

The coordinates for the observations are given by

$$\mathbf{G}_{n \times 2} = \mathbf{U} \mathbf{L}^c \quad (2)$$

and the coordinates for the variables are given by

$$\mathbf{H}'_{2 \times k} = \mathbf{L}^{(1-c)} \mathbf{V}' \quad (3)$$

Biplot-Types are defined by choosing the value for c .

Biplot-Types

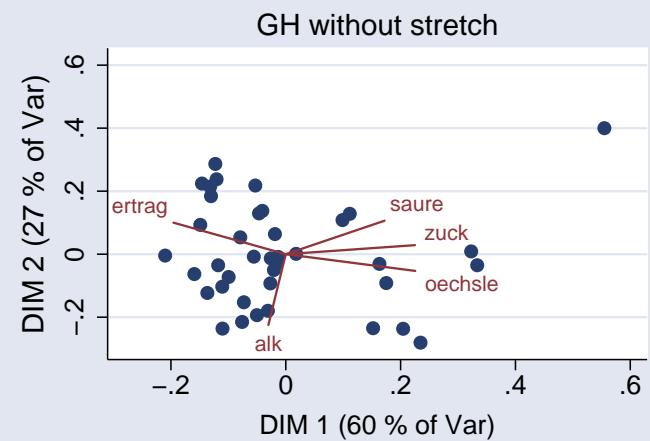
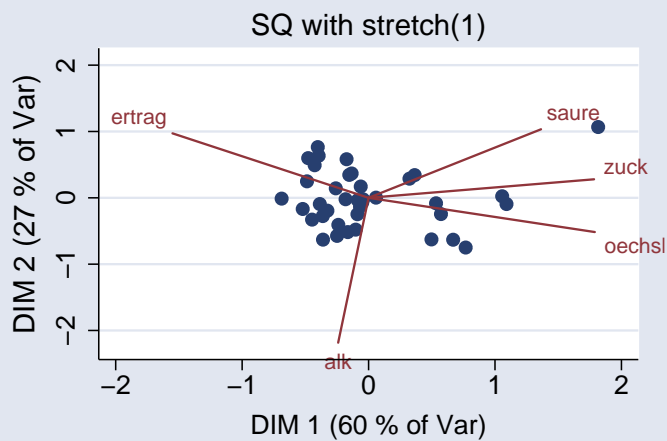
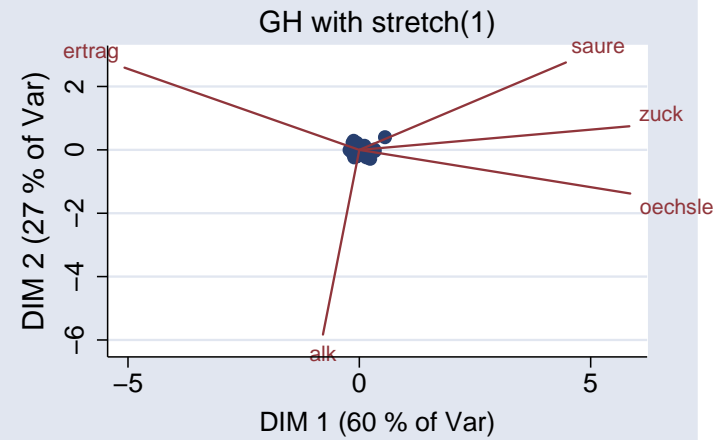
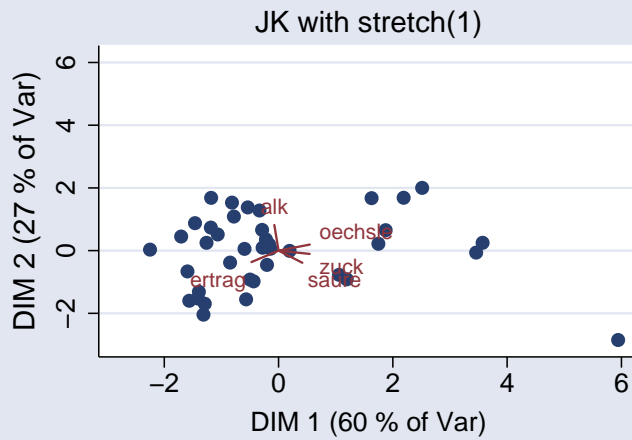
- ⑥ GH-Biplot: $c = 0$
- ⑥ JK-Biplot: $c = 1$
- ⑥ SQ-Biplot: $c = .5$

Note: For $c = 1$ the coordinates for the observations correspond to the first two principal components, and the coordinates for the variables correspond to the first two *Eigenvec-tors*. Therefore `bipplot` calculates a PCA to produce the JK-Biplot.

Biplot-Types

SQ-Biplots are sometimes called *symmetric biplots*. In this type the coordinates of variables and observations tend to be more similar than in the two other types. Regardless of the Biplot-Type, `biplot` automatically chooses a stretch factor for the variable-coordinates making SQ biplots more or less unnecessary.

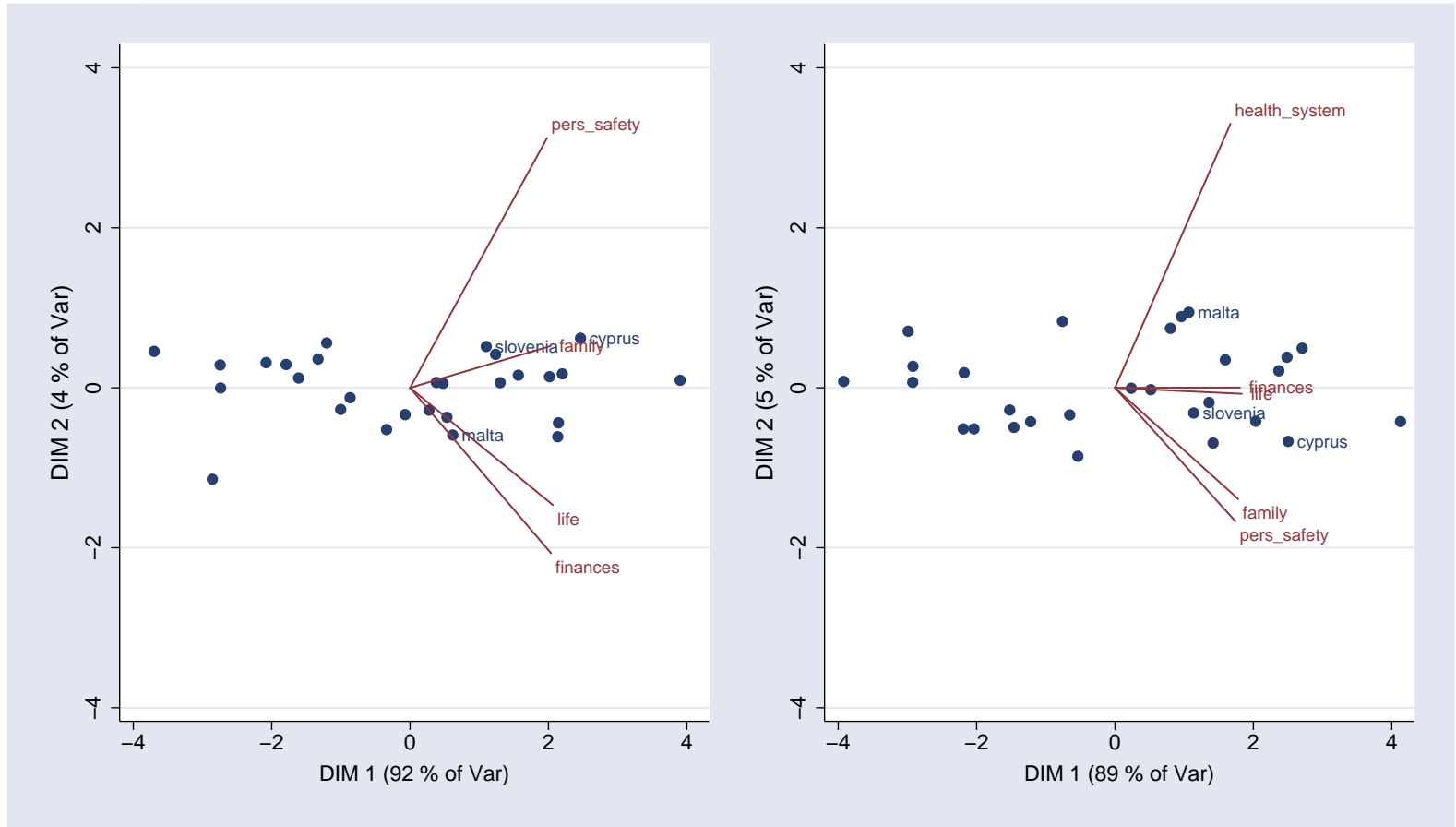
Biplot-Types



Biplot-Types

JK-Biplots are *row metric preserving*, that is, the distances between the objects are more closely approximated in the JK-Biplot than in the other types.

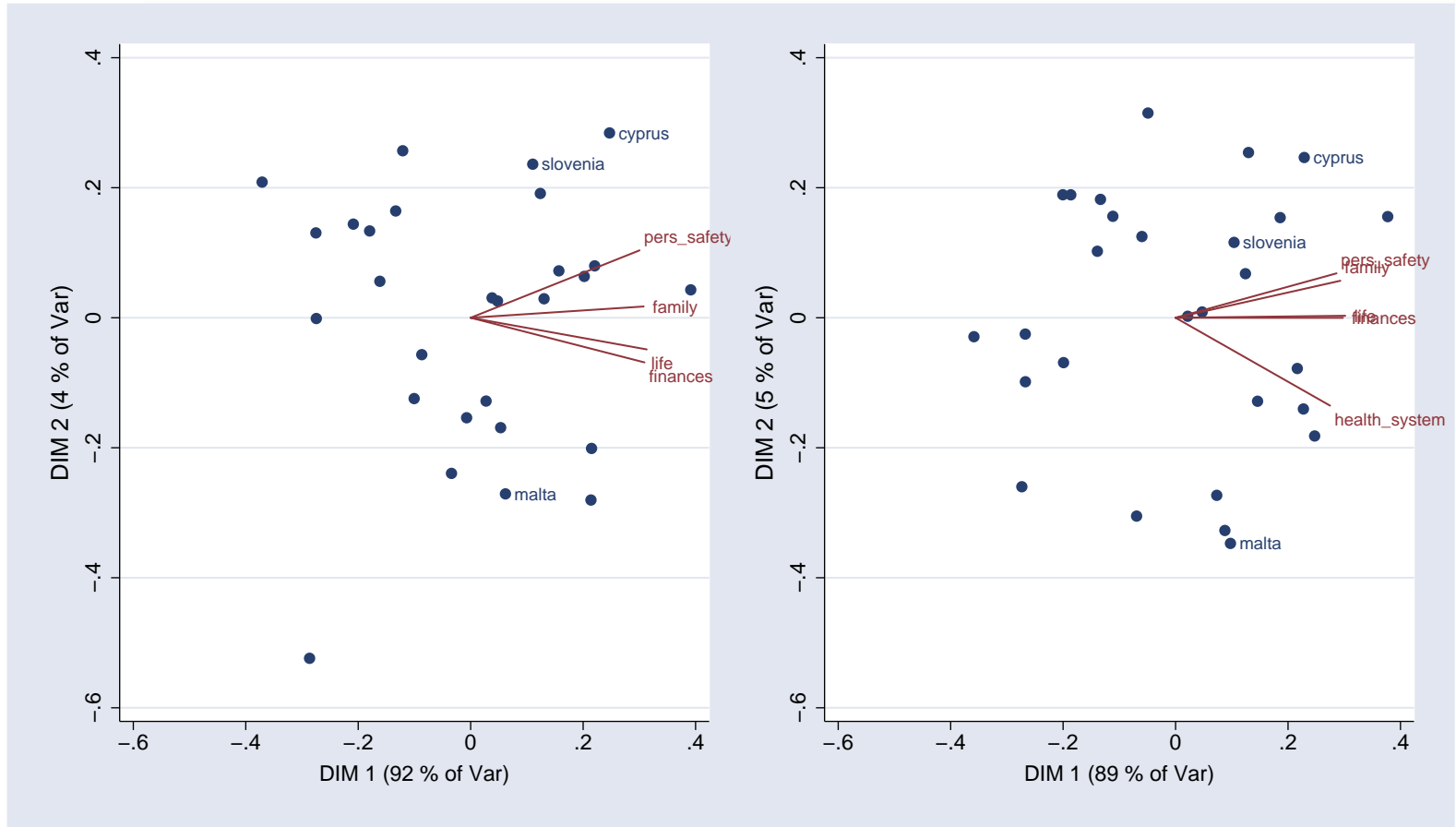
Biplot-Types



Biplot-Types

GH-Biplots are *column metric preserving*, that is, the correlations between the variables are more closely approximated in the GH-Biplot than in the other types.

Biplot-Types



Two new Options

- ⑥ `rv` is used to produce relative variation diagrams. Relative variation diagrams are Biplots for compositional data and compositional data are data sets with constant row-sums and only positive value (like, for example the row percentages of two-way frequency tables). To get a relative variation diagram the data matrix needs to be transformed before producing the Biplot, `biplot` does this transformation for you if you specify `rv`.
- ⑥ `mahalanobis` can be used for GH-Biplots to rescale the graph in a way that the distances between the observations approximates the Mahalanobis distances.