

# Homogeneity Analysis of Durant Bend Sherds

*Jan de Leeuw and C. Roger Nance*

*Version February 17, 2018*

## Abstract

This paper is a non-technical and mainly graphical introduction to *Homogeneity Analysis*, also known as *Multiple Correspondence Analysis*, with an example from archeology.

## Contents

<b>1</b>	<b>Data</b>	<b>1</b>
<b>2</b>	<b>Homogeneity Analysis</b>	<b>2</b>
<b>3</b>	<b>Reciprocal Averaging</b>	<b>4</b>
<b>4</b>	<b>Specifics</b>	<b>8</b>
4.1	Equality Restrictions . . . . .	8
4.2	Binary Variables . . . . .	13
<b>5</b>	<b>Analysis Durant Bend Data</b>	<b>14</b>
	<b>References</b>	<b>15</b>

## 1 Data

The data for this paper are 5888 pot sherds excavated in 1971 and 1975 from several sites in the vicinity of Durant Bend, Dalla County, Alabama. Each sherd was labeled by site and by the depth level in the excavation. In addition each sherd was classified using three binary variables: design (check stamped vs plain), paste color (dark vs light), and thickness (thick vs thin). For the details we refer to Nance (1976) and Nance and De Leeuw (2018). In the following table the sherds are aggregated over sites/depths.

##	CS	Plain	Dark	Light	Thin	Thick
## Ds73S1	118	555	425	248	157	516
## Ds73S2	171	302	304	169	126	347
## Ds73S3	83	156	163	76	67	172
## Ds73S4	36	51	59	28	24	63
## Ds73S5+	25	50	46	29	19	56

## Ds73NAM	203	964	656	511	183	984
## Ds73NUM	196	389	342	243	100	485
## Ds73NLM	164	199	229	134	64	299
## Ds73NBM	74	85	102	57	27	132
## Ds791	11	292	170	133	125	178
## Ds792	17	247	163	101	116	148
## Ds793	24	163	100	87	62	125
## Ds794+	3	39	25	17	18	24
## Au1131	7	34	30	11	30	11
## Au1132	14	40	38	16	32	22
## Au1133+	15	18	25	8	22	11
## Ds98PZ	20	219	166	73	183	56
## DBUCC	11	110	36	85	85	36
## DBLCC	9	45	32	22	47	7
## DBBCC1	13	20	19	14	29	4
## DBBCC2	28	34	45	17	49	13
## DBBCC3	90	79	132	37	147	22
## DBBCC4	83	67	106	44	111	39
## DBBCC5+	66	45	81	30	81	30
## DS971	1	75	60	16	49	27
## DS972	7	67	47	27	51	23
## DS973+	9	45	43	11	38	16

## 2 Homogeneity Analysis

The technique we will use to analyse the Durent Bend data is Homogeneity Analysis (Gifi (1990)), which is more widely known as *Multiple Correspondence Analysis* (Greenacre (1984), Greenacre and Blasius (2006)). We give a graphical introduction to Homogeneity Analysis, without using formulas.

Suppose we have  $m$  categorical *variables*, and that variable  $j$  has  $k_j$  *categories* ( $j = 1, \dots, m$ ). The  $m$  variables *categorize* or *measure*  $n$  *objects*. Variable  $j$  partitions the set of  $n$  objects into  $k_j$  subsets, one subset for each category. Before we get to the analysis of the Durant Bend sherds, we will illustrate the main concepts of our approach with a small example in which three variables partition ten objects. The first two variables have three categories, the last one has two categories.

##	first	second	third
## 01	a	p	u
## 02	b	q	v
## 03	a	r	v
## 04	a	p	u
## 05	b	p	v
## 06	c	p	v

## 07	a	p	u
## 08	a	p	v
## 09	c	p	v
## 10	a	p	v

In Homogeneity Analysis we aim to make a *joint plot* of the objects and the categories of the variables. Joint plots are also known as *biplots* (Gower and Hand (1996)). In a joint plot both objects and categories are represented as points in a low-dimensional space, usually the plane, in such a way that the relations in the data are represented as precisely as possible in the plot. We will specify what we mean by “as precisely as possible” below. Homogeneity Analysis is defined by defining a measure of the loss of information in a certain way, and then choosing the representation that minimizes that loss.

The  $n$  points in the plan, or more generally in  $p$ -dimensional space, representing the objects (sherds) are collected in a matrix of *object scores*. The  $k_j$  points representing the categories of variable  $j$  are in a matrix of *category quantifications* for variable  $j$ . For each variable we can make a *graph plot* in which each of the  $n$  object scores is connected by a straight line to the quantification of the category that this object falls in. Thus there is one line departing from each object point, while the number of lines arriving at a category point is equal to the number of objects in the category. One graph plot has  $n$  lines, all graph plots together have  $n \times m$  lines.

If all these lines have length zero, then all objects coincide with “their” categories for that variable, and thus we have reproduced the data exactly. If there is more than one variable, however, we cannot expect to have such a perfect representation, because objects which are together in a category for one variable may not be together for another variable.

Homogeneity Analysis is defined as the technique that produces a joint plot of objects and category quantification in such a way that the total length of all  $n \times m$  lines in the  $m$  graph plots is as small as possible. Some qualifications are needed, however. We actually minimize the sum of the squared length of the lines, for the same reason that we use the squares of the residuals in a regression analysis. It simplifies the mathematics and the computation to use squared distances. Secondly, we could trivially gain our objective of minimizing line length by collapsing all object scores and category quantifications in a single point, which makes our loss function equal to zero, but is useless in representing or reproducing the data. Thus we need some form of *normalization* to prevent this trivial solution from happening. In Homogeneity Analysis we require the columns of the object score matrix add up to zero, have sum of squares equal to one, and are uncorrelated.

Let’s illustrate this with our small example. We start with a completely arbitrary *initial configuration*. The ten objects are placed at equal distances on a circle, and the categories for each of the variables are on the horizontal axis. This leads to the first three graph plots, which we have superimposed to get the fourth plot with  $n \times m = 30$  lines.

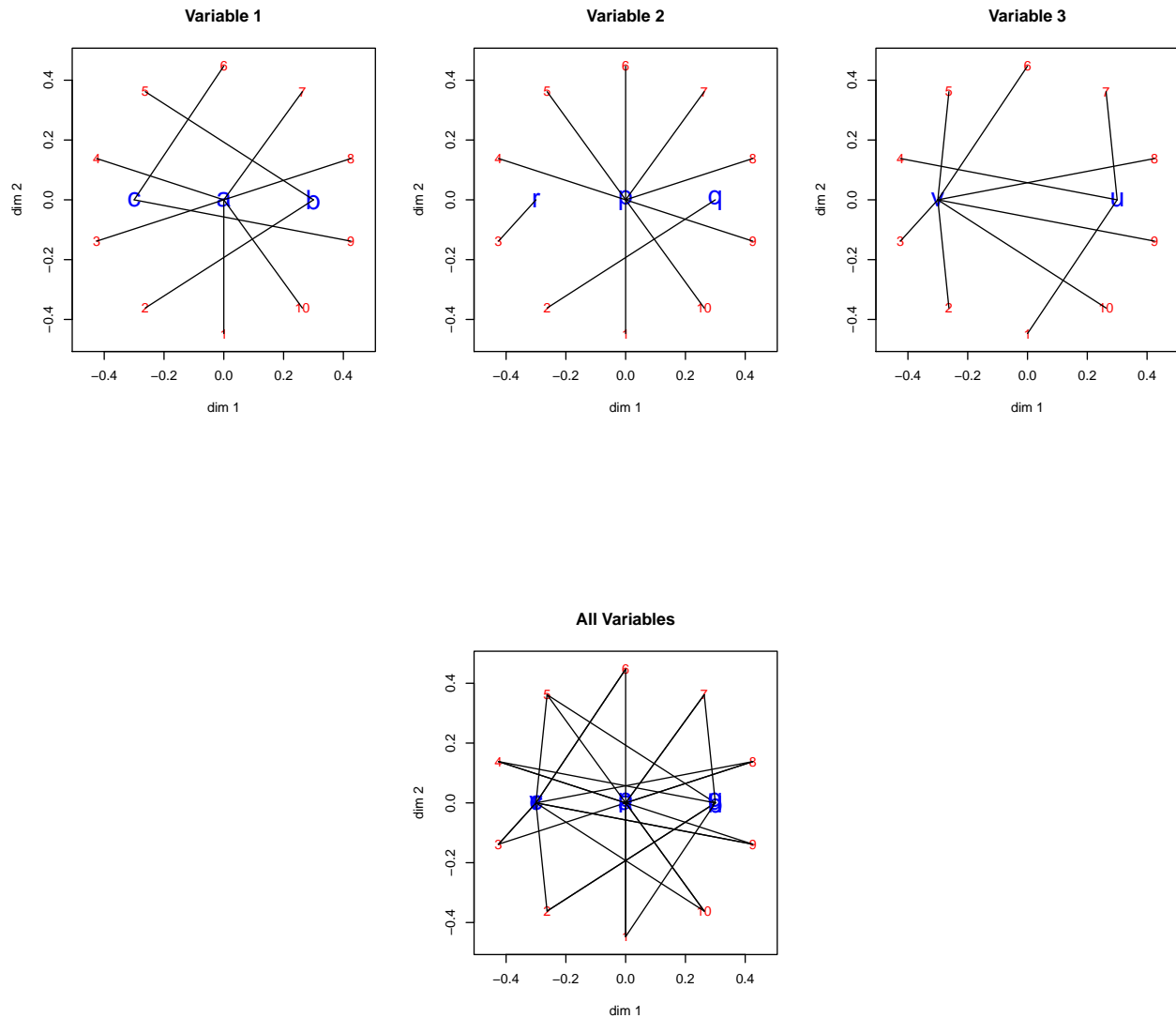


Figure 1: Graph Plots Initial Configuration, Small Example

For this arbitrary initial configuration the loss, i.e. the sum of squares of the line lengths, or the sum of the squared distances between objects and the categories they fall in, is equal to 8.1081098.

### 3 Reciprocal Averaging

In Homogeneity Analysis we minimize loss by what is known as *reciprocal averaging* or *alternating least squares*. We alternate two substeps. The first substep improves the category quantifications for a given set of object scores, the second substep improves and normalizes the object scores for a given set of category quantifications, namely those we have just computed in the first substep. Taken together these two substeps are an *iteration*. So each iteration starts with object scores and category quantifications and uses its two substeps to improve both. Each of the two substeps decreases the loss functions, i.e. the total squared length of

the lines in the graph plots.

The two substeps are both very simple. Let's look at the first one. We compute optimal category quantifications for given object scores by taking the averages (or *centroids*) of the objects scores in each of the categories. The corresponding graph plots are

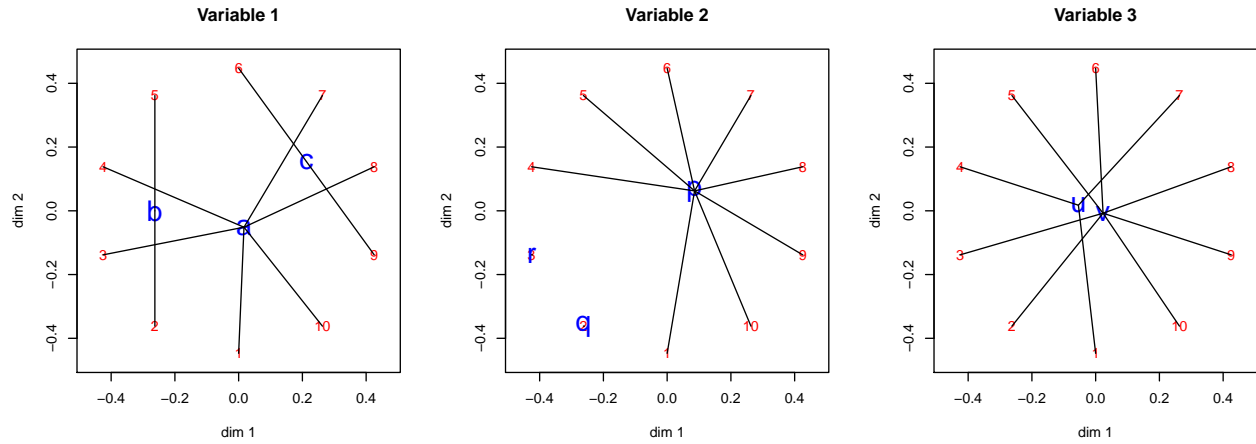


Figure 2: Graph Plots, Iteration 1, substep 1, Small Example

and the loss has decreased to 5.2016654. Note that we have not improved the object scores yet, so they are still their initial configuration, equally spaced on a circle. Also note, in variable 2 for instance, that category quantifications coincide with object scores, and thus contribute zero to the loss, if the object is the only observation in the category. In addition, because category quantifications are averages of objects points, they are in the convex hull of the object points, which means in this figure that they are within the circle. Averaging objects points makes the category quantifications move closer to the origin.

The second substep improves the object scores, while keeping the category quantifications in the locations we have just computed in the first substep. The second substep has itself two substeps, say 2A and 2B. In the substep 2A the score of an object for given category quantifications is computed as the average or centroid of the  $m$  category quantifications the object is in.

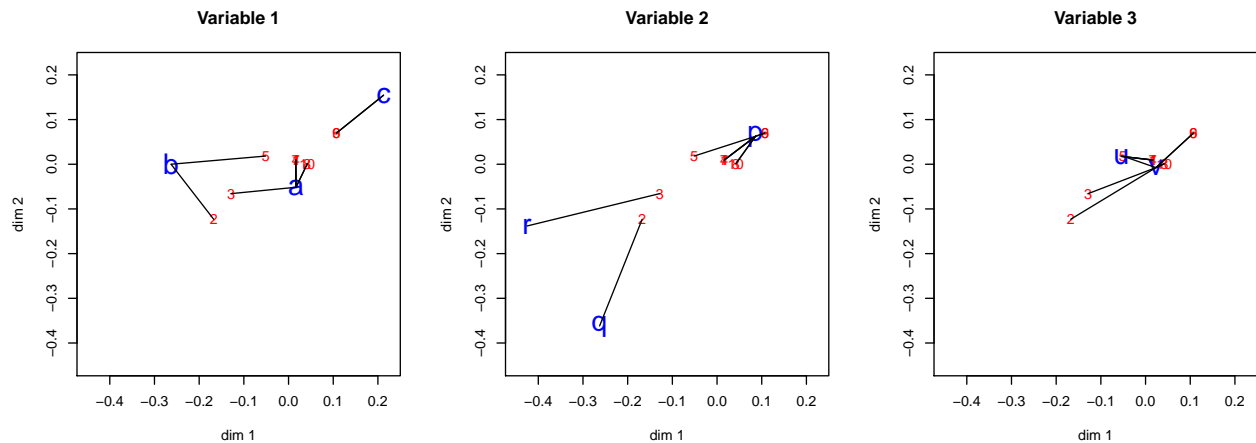


Figure 3: Graph Plots, Iteration 1, substep 2A, Small Example

The loss function is down all the way to 0.4852424. This is not a proper loss value, however, because the object scores are no longer centered, standardized, and uncorrelated, and that was a Homogeneity Analysis requirement. Substep 1 shrinks the object scores towards the origin by averaging, substep 2A takes the resulting category quantifications and shrinks them more by even more averaging. Thus in substep 2B we have to renormalize the object scores such that they are centered, standardized, and uncorrelated. This gives

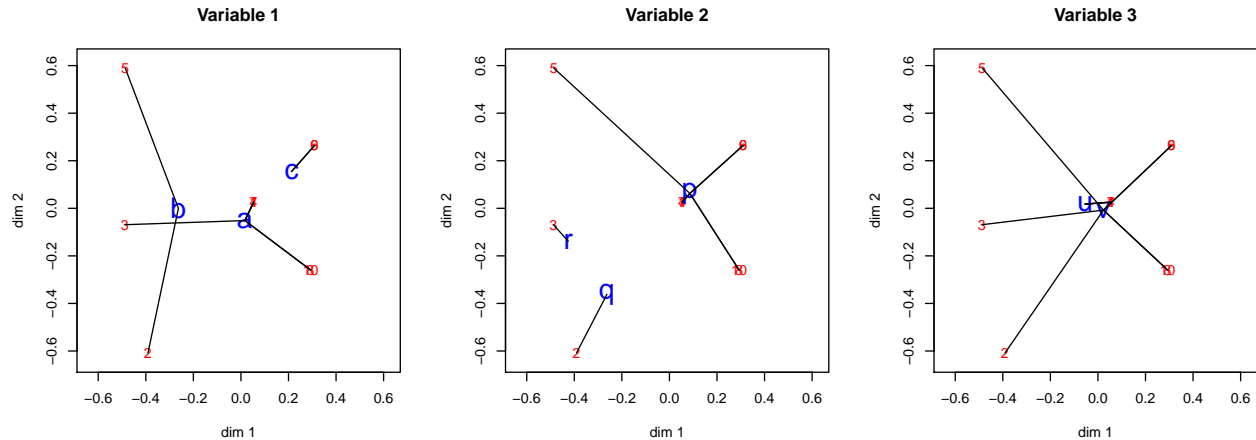


Figure 4: Graph Plots, Iteration 1, substep 2B, Small Example

Loss, which is now the proper loss for a normalized configuration, has decreased to 4.5538169.

Now that we have new category quantifications and new suitably normalized object scores we can start the next iteration, and again improve both in two substeps. Ultimately, after a certain number of iterations, there is no change any more from one iteration to another, and we have reached the optimal solution. In other words, there is convergence, and our Homogeneity Analysis is finished. Note that the renormalization in step 2B is necessary, because without it both object scores and category quantifications would become smaller and smaller, and converge to the origin. Of course the origin does have loss zero, but it is never a proper description of the data.

The optimal graph plots, after the iterations have converged, are

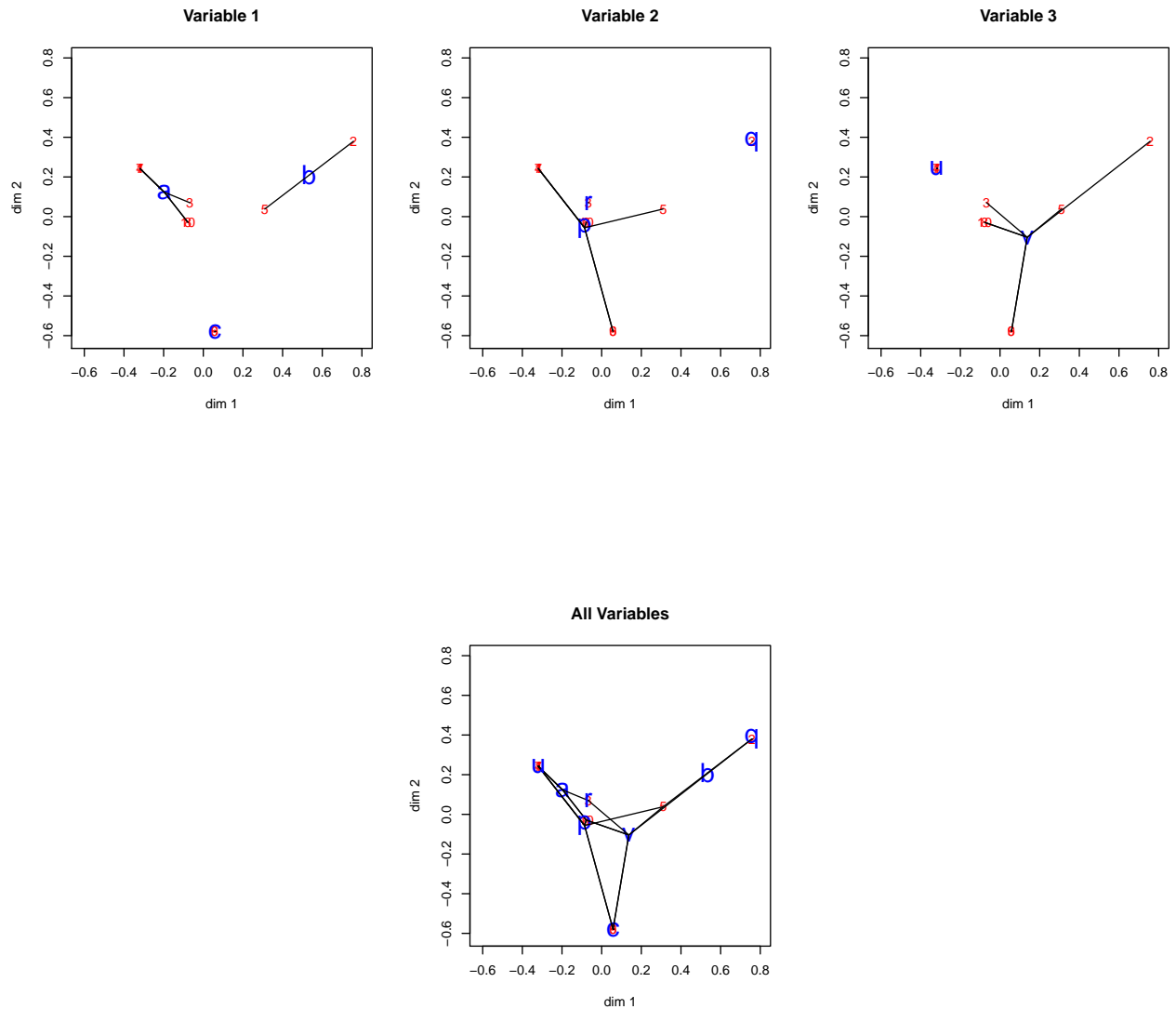


Figure 5: Graph Plots, Optimum Solution, Small Example

The minimum loss for these data is 2.8377218.

Because the object scores are in deviations from the mean, and the category quantifications are weighted means of object scores it follows that category quantifications are in deviations from the weighted mean, with weights equal to the marginals of the variable. Thus category quantifications for each variable are distributed around the origin.

In Homogeneity Analysis the graph plots for individual variables are often called *star plots*, because the optimal category quantification is the centroid of the scores of the objects in the category. Thus it is somewhere in the middle of a bunch of objects, which are connected to it by lines. Thus the subset of the graph for each category is a star graph, and the corresponding plot for the variable with these stars is a star plot. The top three plots in figure 5 are examples of such star plots. One could formulate the objective of Homeogeneity Analysis as finding normalized object scores in such a way that the stars (over all categories of all variables) are as small as possible. Or, in yet another formulation, we want to maximize

the between-category variation and minimize the within-category variation.

Variables with a small star, which is necessarily close to the origin, have poor discrimination power. The average object score for each category is about the same. In general, categories with a large number of observations will have an average close to the average of all observations, and thus they will be close to the origin. And, conversely, categories with a small number of observations will tend to be relatively far from the origin.

## 4 Specifics

### 4.1 Equality Restrictions

Now, going back to the Durent Bend data, we do not have the values of all 5888 sherds on the three variables. The sherds are aggregated over various site/depth combinations and the original data cannot be recovered from the aggregated table. But the framework of Homogeneity Analysis can still be applied by using *equality restrictions* (Van Buuren and De Leeuw (1992)). The only thing added is that we require that sherds in the same site/depth get the same object score. Or, geometrically, all sherds in the same site/depth are mapped into the same point in the joint plot.

The Homogeneity Analysis loss function is still minimized in two steps. The first step, updating the category quantifications, is still the same as in Homogeneity Analysis without equality restrictions. The second step, which updates the object scores, now has three substeps instead of two. In substep *2A* we compute the average of the category quantifications of the categories the sherd is in. In substep *2B* we replace these tentative object scores for the sherds by the site/depth averages, and in substep *2C* we normalize the object scores, making them centered, standardized, and uncorrelated.

The graph plots on unconstrained Homogeneity Analysis must now be replaced by plots of *valued graphs*. For any variable each site/depth point is now connected to all category points, and the edge connecting the object and category point has a value equal to the number of sherds in the category.

As an example we use the GALO data, which has been used innumerable times before as a Homogeneity Analysis example, and is in the `Gifi` package in R (P. Mair and De Leeuw (2017)).

```
galo Gifi    R Documentation
GALO dataset
```

#### Description

The objects (individuals) are 1290 school children in the sixth grade of elementary scho

#### Usage



galo  
Format

Data frame with the five variables Gender, IQ, Advice, SES and School. IQ (original range)

SES:

LoWC = Lower white collar; MidWC = Middle white collar; Prof = Professional, Managers; S

Advice:

Agr = Agricultural; Ext = Extended primary education; Gen = General; Grls = Secondary sc

## References

Peschar, J.L. (1975). *School, Milieu, Beroep*. Groningen: Tjeek Willink.

We first ignore the school variable, and analyze the four variables Gender, IQ, Advice, and SES. Separate joint plots for the four variables, with both object scores and category qualifications, are in figure 6. We do not make star plots (by drawing the lines from the object points to the category points they are in) in this case, because 1290 lines in a plot just create a big black blob. The joint plots show a curved one-dimensional solution with good students on the left and poor students on the right. Both IQ and Advice differentiate students well (because teachers used IQ in their secondary education advice), which means they will have the smallest stars. Girls are better students than boys, and SES mainly contrasts the two extremes categories PROF and UNSK.

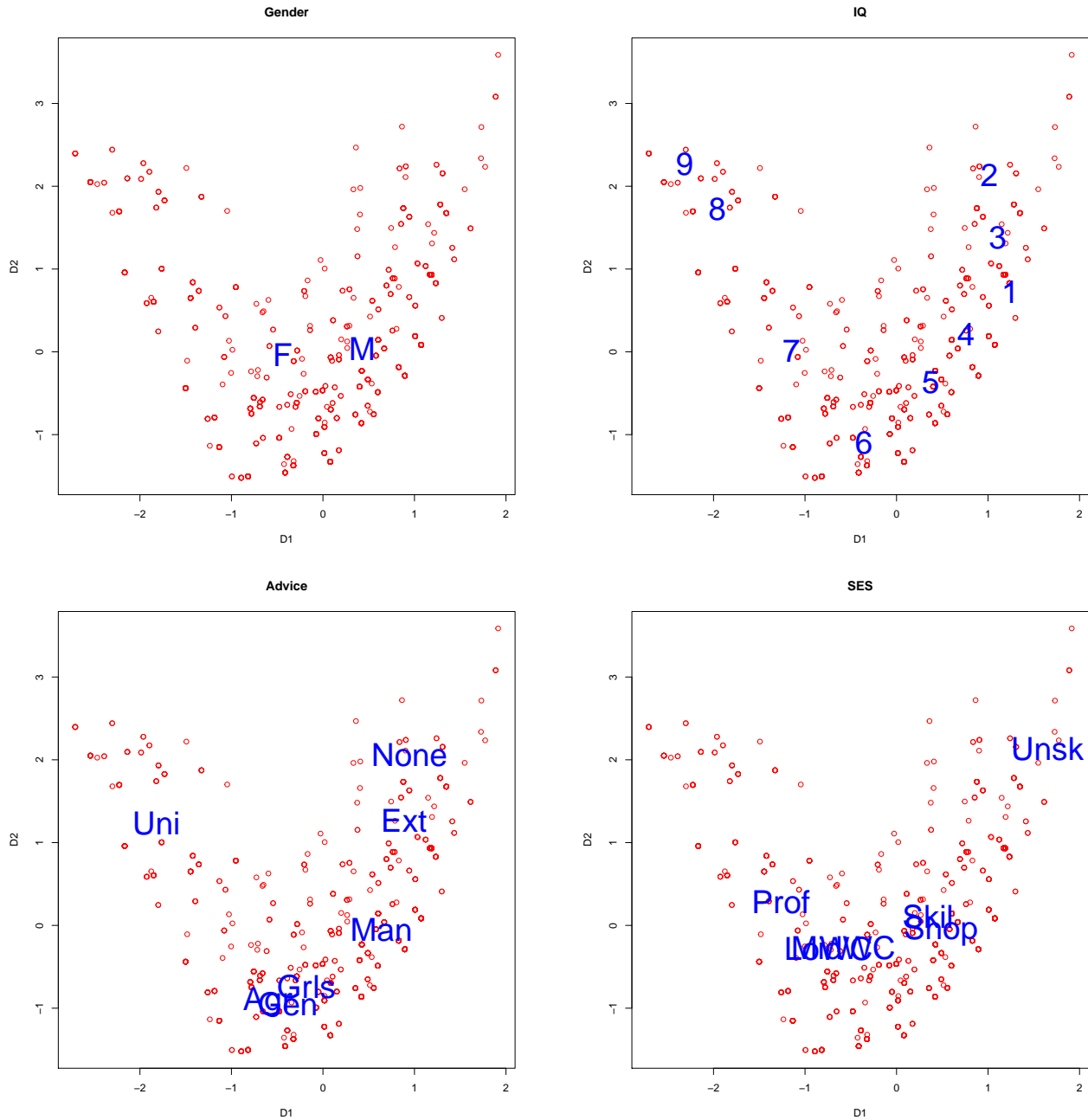


Figure 6: Joint Plots, GALO Example

The analysis does not use the school variable at all. In the terminology of Gifi “school” is a *passive variable*. We now repeat the analysis, requiring that students in the same school get the same object score. Computationally this is easiest to do using the the R package **anacor** (De Leeuw and Mair (2009)). In this analysis each school get an object score, and these scores are plotted in figure 7.

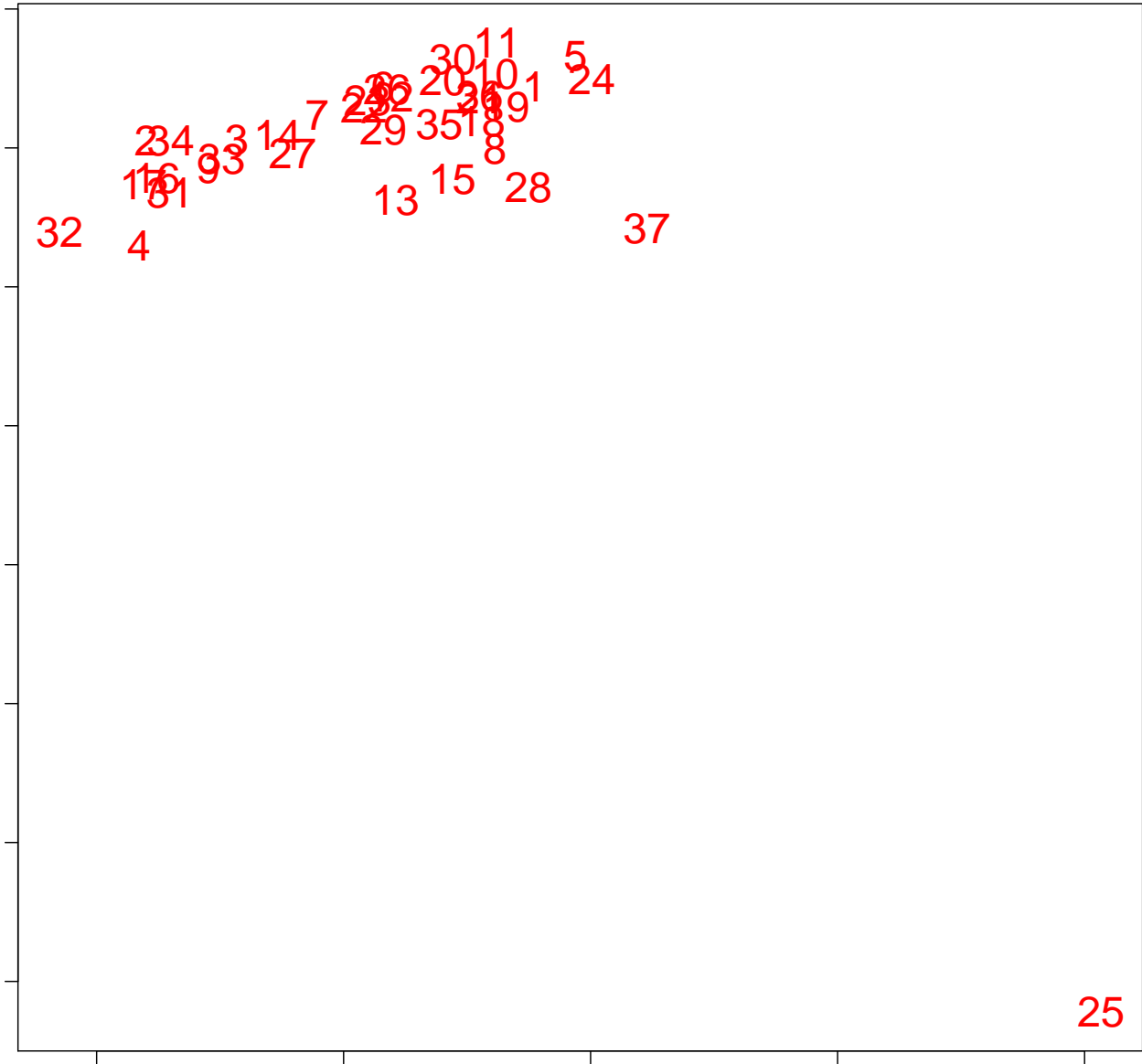


Figure 7: Object Score Plot, School Constraints, GALO Example

School 25 is an outlier, but otherwise schools are dispersed in ways that can possibly be interpreted.

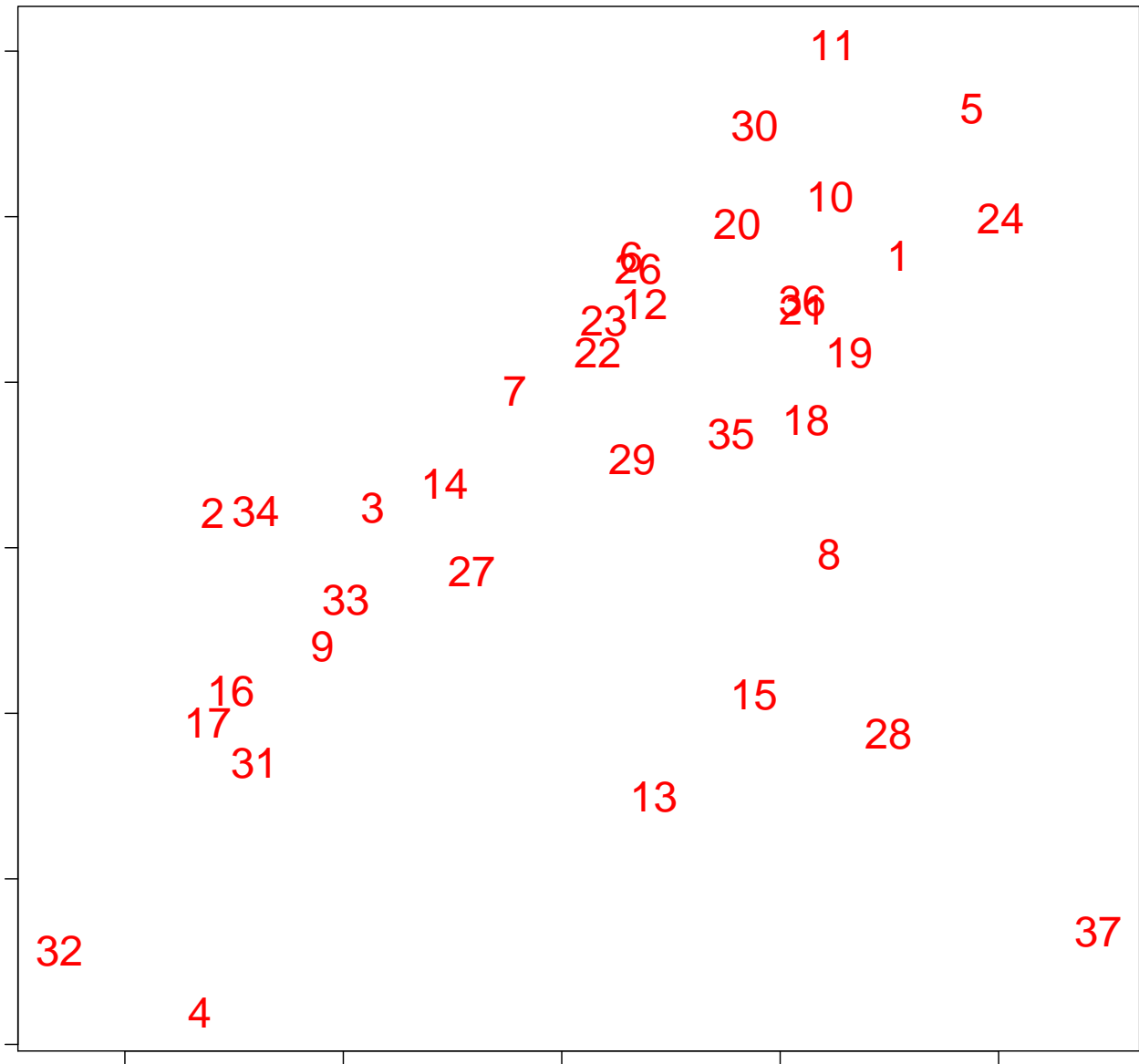
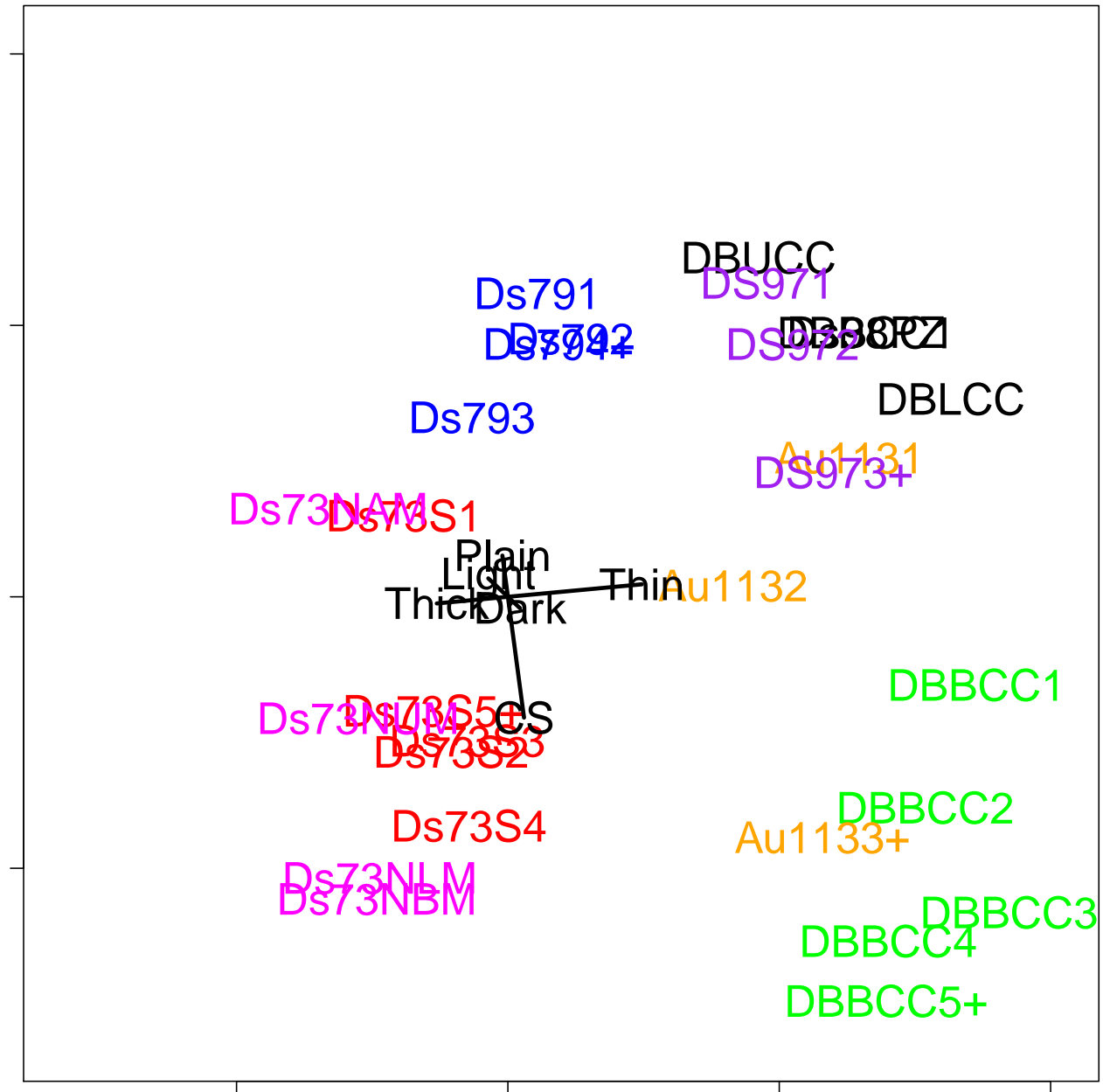


Figure 8: Object Score Plot, School Constraints, Without School 25, GALO Example

There is another less radical way to incorporate “school” into the analysis, by using it as what is commonly known as a *supplementary variable*. Such a variable does not enter into the analysis, but we can still compute the category quantifications as centroids of object scores in the category. Thus we can make star plots for passive variables that have not been used in the analysis.



## 5 Analysis Durant Bend Data

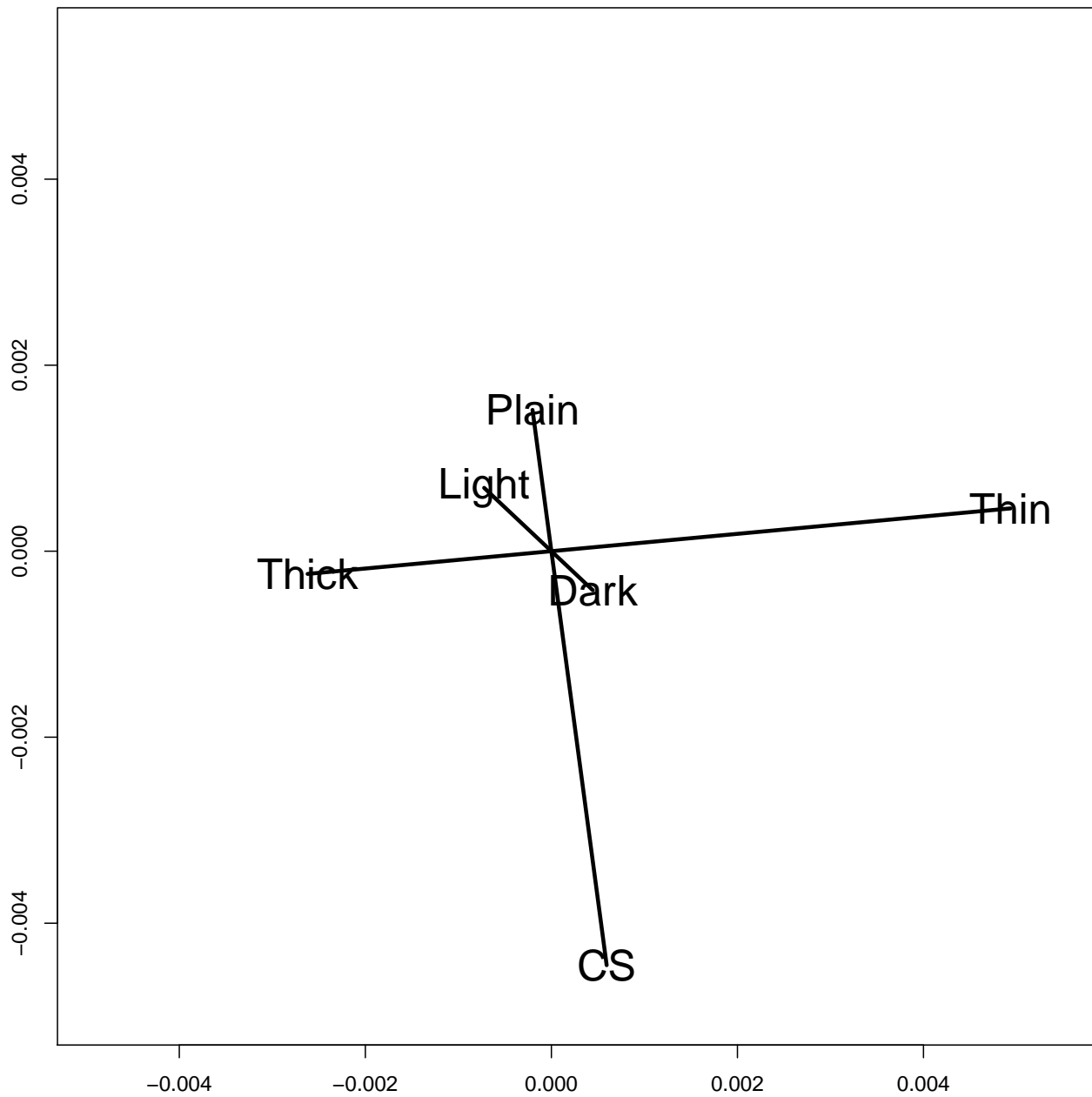


```

par (pty = "s")
y <- rbind(y1, y2, y3)
mA <- 1.1 * max (y)
mI <- 1.1 * min (y)
plot(0, xlim=c(mI,mA), ylim = c (mI, mA), type = "n", xlab = "", ylab = "")
lines(y1, lwd = 3)
text(y1, row.names(y1), cex = 2)
lines(y2, lwd = 3)
text(y2, row.names(y2), cex = 2)

```

```
lines(y3, lwd = 3)
text(y3, row.names(y3), cex = 2)
```



## References

De Leeuw, J., and P. Mair. 2009. "Simple and Canonical Correspondence Analysis Using the R Package Anacor." *Journal of Statistical Software* 31 (5): 1–18. <http://www.stat.ucla.edu/>

~deleeuw/janspubs/2009/articles/deleeuw\_mair\_A\_09b.pdf.

Gifi, A. 1990. *Nonlinear Multivariate Analysis*. New York, N.Y.: Wiley.

Gower, J.C., and D.J. Hand. 1996. *Biplots*. Monographs on Statistics and Applied Probability 54. Chapman; Hall.

Greenacre, M.J. 1984. *Theory and Applications of Correspondence Analysis*. New York, New York: Academic Press.

Greenacre, M.J., and J. Blasius, eds. 2006. *Multiple Correspondence Analysis and Related Methods*. Chapman; Hall.

Mair, P, and J. De Leeuw. 2017. *Gifi: Multivariate Analysis with Optimal Scaling*. {<https://R-Forge.R-project.org/projects/psychor/>}.

Nance, C. Roger. 1976. “The Archeological Sequence at Durant Bend, Dallas County, Alabama.” Special Publication 2. Orange Beach: Alabama Archeological Society.

Nance, C. Roger, and J. De Leeuw. 2018. “Statistical Perspectives on Woodland Cultures in Central Alabama.”

Van Buuren, S., and J. De Leeuw. 1992. “Equality Constraints in Multiple Correspondence Analysis.” *Multivariate Behavioral Research* 27: 567–83. [http://www.stat.ucla.edu/~deleeuw/janspubs/1992/articles/vanbuuren\\_deleeuw\\_A\\_92.pdf](http://www.stat.ucla.edu/~deleeuw/janspubs/1992/articles/vanbuuren_deleeuw_A_92.pdf).