

# Mosaic Mind Games — Visualising Categorical Data

#### Antony Unwin

Augsburg University, Germany

Antony Unwin

Mosaic Mind Games

Auckland University 11th April 2013





Auckland University 11th April 2013

### Mosaicplots



- Variable category combinations are represented by rectangles. There are gaps between rectangles (ideally smaller by level).
- Rectangle area is (almost always) proportional to frequency. Rectangles may have equal width (height), so that height (width) is proportional to frequency.
- Rectangles may be aligned in various ways, rotated, coloured.
- Mosaicplots need to be interactive.

Mosaic Mind Games

Antony Unwin Auckland University 11th April 2013

# Aside: Titanics in R



- Titanic {*datasets*}, Titanic {*effects*}
- TitanicMat {*RelativeRisk*}
- titanic {*alr4*}, titanic {*prLogistic*}, titanic {*msme*}
- titanic.dat {*exaxtLoglinTest*}
- titan.Dat {*elrm*}
- etitanic {*earth*}
- ptitanic {*rpart.plot*}
- Lifeboats {*vcd*}

and maybe there are more I have not found ....

Mosaic Mind Games



Auckland University 11th April 2013

So take care!

### Titanic dataset



2201 passengers and crew by gender, age (child/adult), ship's class (1st, 2nd, 3rd, crew), survived or died.

(R. J. MacG. Dawson, J. Statistics Education <u>3</u> no 3, 1995)













The *really* interesting thing about the neural network's solution to the problem isn't that it reached any kind of useful error rate (which it didn't), but that its weights encode the relative importances of the various inputs.







The configuration clearly distinguishes two groups, the outer points representing children and the people scattered around the origin representing adults (upper right panel of Figure 2.4). One reason for this result is that there were not many children on board of the ship, and objects of low frequency tend to be located in the periphery of the plot.

The lower right panel of Figure 2.4 shows that most survivors are located on the left-hand side of the first axis. Since all women (lower left panel of Figure 2.4) as well as most children (upper right exiet of Figure 2.4) are located in this very same area, the PIONEER analysis shows that the rule women and children first' seems to have been applied in the rescue operation of the sinking Titanie. Fig. By, there is a slight indication that first and second class passengers were rescued more often then the others (upper left panel of Figure 2.4).

Figure 2.4. Labeled profile scores plots of PIONEER analysis of Titanic survival data.

Patrick J.F. Groenen, Jacques J.F. Commandeur Department of Data Theory, Leiden University





Figure 9 shows the frequencies of the background variables, Class, Gender and Age by the sizes of the boxes. It also shows the association between Age and Class–Gender combinations by shading. There were no children among the crew, and the overall proportion of children was quite small (about 5 %). But among the passengers, the proportion of children increases from first class to third class. The large positive residuals for children among the 3rd class passengers likely represents families traveling or emmigrating together.

Figure 10 shows an initial four-way mosaic for the full table, and fits the model which asserts that survival is independent of Class. Gender and Aga jointly. This is the minimal and model when the first three variables are explanatory. It is clear that greater proportions of women survived than men in all classe, but with greater proportions of women survived than men in all classe, but with greater proportions of women survived than men in all classe, but with greater proportions of women survived than men in all classe, but with greater proportions of women survived than men in all classe, but with greater proportions of women survived that men appendent of the survived survived by the survived by the survived two increases with economic class. However, this model for very poorly ( $G^2(15) \equiv 6/1.96$ ), and we may try to fit a more adequate model by adding associations between survival and the explanatory vari- ables.













95% CI

Lower Upper

1.83 4.51 0.15 0.01 0.04

0.07 1.68 0.13 2.19 1.65 21.21

5.26 1.10 25.26 0.26 3.77

# Titanic models (Simonoff)

Table of models						Gender*Class + Age model					
Model G E A E, G A, G E, A E, A E, A, G E, G, EG A, G, AG E, A, EA	G 434.5 180.9 19.6 540.5 440.4 206.5 559.4 605.7 456.7 235.7	Predictors 1 3 1 4 2 4 5 7 3 6	D .40 .28 .05 .49 .42 .30 .52 .50 .43 .30	H-L p .000 .000 .001 .000	$\begin{array}{c} AIC_{C} \\ -430.5 \\ -172.9 \\ -15.6 \\ -530.6 \\ -434.9 \\ -196.5 \\ -547.4 \\ -589.7 \\ -448.7 \\ -221.7 \end{array}$	Logistic Regressio Predictor Constant Economic First class Second class Third class Age grou Child Gender	n Table Coef 1.8971 1.6608 -0.0199 -2.2247 1.0537	StDev 0.6191 0.8003 0.6869 0.6370 0.2304	Z P 3.06 0.002 2.08 0.038 -0.03 0.977 -3.49 0.000 4.57 0.000	Odds Ratio 5.26 0.98 0.11 2.87	95% ( Lower 1.10 0.26 0.03 1.83
E. A. G. EG E. A. G. AG E. A. G. AG E. A. G. EA, EG E. A. G. EA, EG E. A. G. EA, AG E. A. G. EA, AG E. A. G. EA, EG, AG	626.1 595.1 577.4 670.3 634.7 606.9 672.0	8 7 6 10 9 8 11	.53 .52 .53 .53 .53 .54	.941 .000 .944 .992 .000 1.000	-608.0 -579.0 -563.4 -648.2 -614.6 -588.9 -647.9	Male Economic*Gender First class*Male Second class*Male Third class*Male	-3.1469 -1.0862 -0.6379 1.7763	0.6245 0.8197 0.7250 0.6522	-5.04 0.000 -1.33 0.185 -0.88 0.379 2.72 0.006	0.04 0.34 0.53 5.91	0.01 0.07 0.13 1.65
Mosaic Mind Games					Antony Unwin			Auckland University 11th April 2013			



#### Two of the alternatives







Equal binsize plot Titanic survival rates by class gender, age Women to the left, men to the right Doubledecker plot Titanic survival rates by class and gender Women to the left, men to the right

Mosaic Mind Games

Antony Unwin Auckland University 11th April 2013

Categorical data and visualizations thereof

Mosaic Mind Games

Antony Unwin

Auckland University 11th April 2013

# Categorical data



- Nominal data (perhaps grouped, e.g. geographically)
  - Occupation, Experimental treatments, Cities, ...
- Binary or logical data
  - Gender, Yes/No, True/False, ...
- Ordinal data
  - Survey responses, Income group, Fitness, ...
- Discrete data (and discretised continuous data)
- Examples: Titanic, Rochdale, Divorce, Bowling Alone, ...

Mosaic Mind Games

Antony Unwin Auc

Auckland University 11th April 2013

## Plots for categorical data



- barcharts, stacked barcharts, dodged barcharts
- piecharts, agreement plots (Bangdiwala)
- fourfold displays, sieve diagrams, association plots, cpcp
- 3-d and trellis versions
- mosaicplots
- familes of plots in R
  - strucplot (Meyer, Zeileis, Hornik) vcd package
  - productplots (Wickham, Hofmann) productplots package

Mosaic Mind Games

Antony Unwin

Auckland University 11th April 2013







#### 3D mosaicplot (vcdExtra)

For those of a sensitive disposition I would suggest avoiding Figure 13.8, which apparently shows a three-dimensional mosaicplot. Up to this point in the book I had agreed with Paul Murrell's statement in his Preface that "no plot type is all bad".



#### Bangdiwala agreementplot (vcd)

@Antony: Since this is a Fig. I take it that you are prepared to argue that fluctuation diagrams are an alternative (maybe better) display for agreements. I caution you not to say that in print, because you would be wrong... Michael Friendly

Mosaic Mind Games Antony Unwin

Auckland University 11th April 2013

<figure><figure><figure><complex-block><complex-block><table-row><table-row><table-row><table-row><table-row><table-row><table-row><table-row><table-row></table-row>



#### Another dataset?

- Possibilities
  - *wong.df* in James Curran's *dafs*

- nhanes in Thomas Lumley's survey

– *wffc* in Thomas Yee's *VGAM* 



- diabetes, murder, NORC, Auckland .... in Stats 330
- Decided to use *Intergenerational inequality* from StatsChat

   GSS (USA) with education, parents' education, age, sex, family income, survey weighting, ...

Antony Unwin

Mosaic Mind Games

Auckland University 11th April 2013

# Intergenerational inequality



Mosaic Mind Games

Antony Unwin Auckland University 11th April 2013

```
T CONSCIENCE
```





# Classical mosaicplots













# Mosaicplots also include



- Residual plots (by expected/observed, association plots)
- Relative multiple barcharts
- Multiple spineplots
- Treemaps

though no single software package implements all (yet)

Antony Unwin

Mosaic Mind Games

Auckland University 11th April 2013

# Mosaicplot options



- Choice of variables
- Order of variables
- Whether each variable is horizontal or vertical
- Form of mosaicplot
- Orders of categories within nominal variables, ordering direction for ordinal variables
- Grouping categories
- Display options: spacing between levels and between categories, plot size, aspect ratio, colour

Mosaic Mind Games

Antony Unwin Au

Auckland University 11th April 2013

# Numbers of mosaicplots



Titanic

≥8	Variants	8	
2 <sup>m</sup> - 1	Choice from <i>m</i> variables	7	
r!	Orderings of <i>r</i> variables	24	
2 <sup>r</sup>	Directions of variables (horizontal/vertical)	16	
$\prod c_j !$	Orderings of categories within variables	24*2*2*2	
	(or 2 <sup>r</sup> for direction of ordinal variables)		
?	Aggregations of categories		
?	New derived variables		
Mosaic Mind Ga	ames Antony Unwin Auckland University 11th	n April 2013	



#### Variant choices



- Classical mosaicplots: for cumulative rates
- Residual plots: for supporting model building
- Fluctuation diagrams: for sparse structures
- Multiple barcharts: for non-binary target variables
- Same binsize: for rates across all groups and missing values
- Doubledecker plots: for rates across all groups with cell sizes
- Relative multiple barcharts: for distribution shape

Antony Unwin

• Treemaps: for splitting by different variables

Mosaic Mind Games

Auckland University 11th April 2013

# Design principles (1)



- Variable ordering
  - Target variable should usually be last
  - Binary target variables are best displayed using linking
  - A grouping variable should be first, possibly rotated
  - Comparisons and context determine the order of conditioning variables (+ overall height/width)
  - Variables with unequal distributions are better early (?)

Mosaic Mind Games	
-------------------	--

Auckland University 11th April 2013

Design principles (2)



- Category ordering
  - can be determined by context
  - by what you want to compare
  - can sort by count, absolute proportions, relative proportions
- Vary aspect ratio of cells
  - square (fluctuation diagrams), otherwise height > width
- One graphic is usually not enough

Mosaic Mind Games

Antony Unwin

Auckland University 11th April 2013



### Mondrian



- Mondrian for interactive graphical analysis
  - one of the Augsburg Impressionists
  - -stats.math.uni-augsburg.de/Mondrian/
  - -for Windows, Unix, MacOS
  - -by Martin Theus



Mosaic Mind Games

Antony Unwin Auckland University 11th April 2013

## Summary



- Categorical data are difficult to visualise
- Several related plots are more effective than one single plot
- Mosaicplots are a general, flexible family of displays for categorical data though
  - they are often puzzling to interpret
  - and they need thoughtful design
  - so sometimes they seem more like mind games