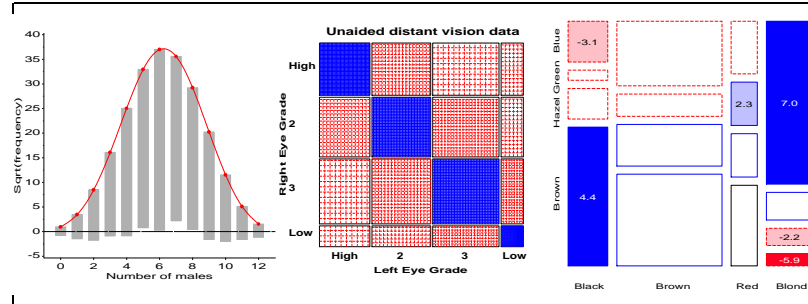


Categorical Data Analysis with Graphics



Michael Friendly

York University, Toronto, <friendly@yorku.ca>



SUGI 28
April, 2003



Categorical Data Analysis

Methods of analysis for categorical data fall into two main categories:

■ Non-parametric, randomization-based methods

- make minimal assumptions
- useful for hypothesis-testing
- SAS: PROC FREQ
 - Pearson Chi-square
 - Fisher's exact test (for small expected frequencies)
 - Mantel-Haenszel tests (ordered categories: test for *linear* association)

■ Model-based methods

- Must assume random sample (possibly stratified)
- Useful for estimation purposes
- Greater flexibility; fitting specialized models (e.g., symmetry)
- More suitable for multi-way tables
- SAS: PROC LOGISTIC, CATMOD, GENMOD, INSIGHT (Fit YX)
 - estimate standard errors, covariances for model parameters
 - confidence intervals for parameters, predicted $Pr\{\text{response}\}$

SUGI 28

2

Michael Friendly

Outline

- Overview: Categorical Data *and* Graphics
- Methods for discrete distributions
 - Hanging rootograms
 - Robust distribution plots
- Methods for two-way frequency tables
 - Fourfold displays
 - Sieve diagrams
- Mosaic displays and loglinear models for n -way tables
 - Mosaic displays
 - Mosaic matrices
- Logistic and logit regression
 - Logit plots, effect plots
 - Diagnostic plots

Color version of these slides:
<http://www.math.yorku.ca/SCS/sugi/sugi28.pdf>

SUGI 28

1

Michael Friendly

Graphical Methods for Categorical Data

If I can't picture it, I can't understand it.

Albert Einstein

Getting information from a table is like extracting sunlight from a cucumber.

Farquhar & Farquhar, 1891

■ Tables vs. Graphs

- Tables are best suited for *look-up*— read off exact numbers
- Graphs are better for showing *patterns, trends, anomalies*, making *comparisons*
- Visual presentation as *communication*: what do you want to say?

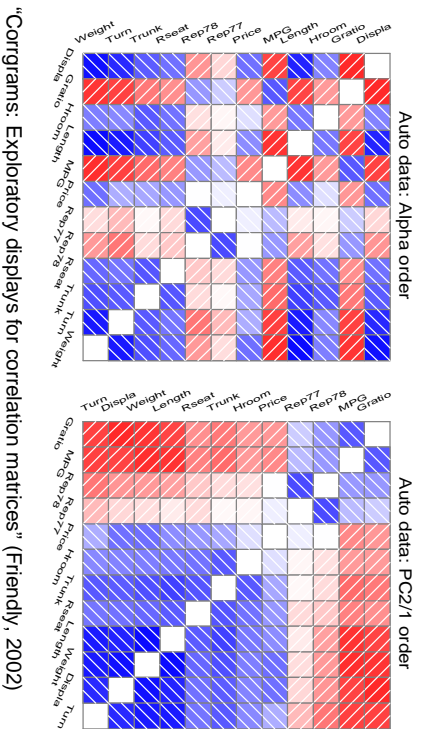
SUGI 28

3

Michael Friendly

■ Principles of Graphical Displays

- Effect ordering (Friendly and Kwan, 2002)— In tables and graphs, sort unordered factors according to the effects you want to see/show.



SUGI 28

4

Michael Friendly

- Effect ordering and high-lighting for tables (Friendly, 2000)

Table 1: Hair color - Eye color data: Effect ordered

Eye color	Hair color		
	Black	Brown	Red
Brown	68	119	26
Hazel	15	54	14
Green	5	29	14
Blue	20	84	17

Table 2: Hair color - Eye color data: Alpha ordered

Eye color	Hair color		
	Blond	Black	Red
Blue	94	20	17
Brown	7	68	26
Green	10	15	14
Hazel	16	5	14

Model: Independence: [Hair][Eye] $\chi^2(9) = 138.29$

Color coding: $n < -4$ $n < -2$ $n < -1$ 0 $n > 1$ $n > 2$ $n > 4$

n in each cell: $n < \text{expected}$ $n > \text{expected}$

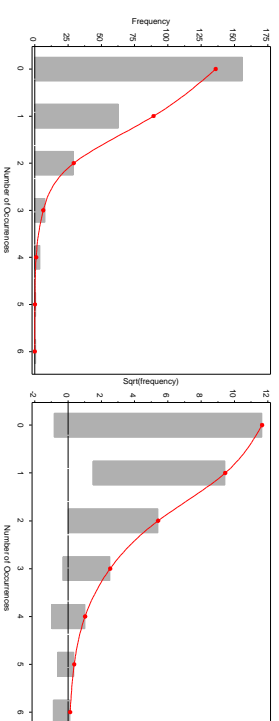
SUGI 28

5

Michael Friendly

■ Comparisons— Make visual comparisons easy

- Visual grouping— connect with lines, make key comparisons contiguous
- Baselines— compare data to model against a line, preferably horizontal



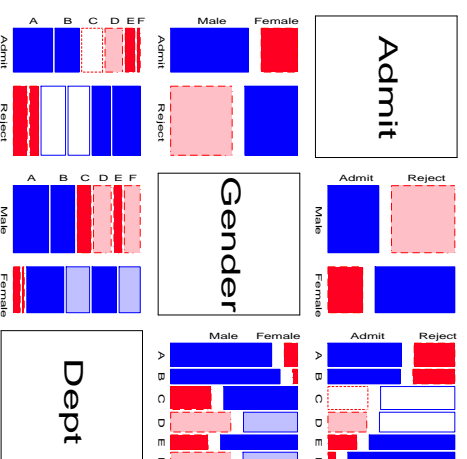
SUGI 28

6

Michael Friendly

■ Comparisons— Make visual comparisons easy

- Small multiples (Tuft, 1983)— combine stratified graphs into coherent displays



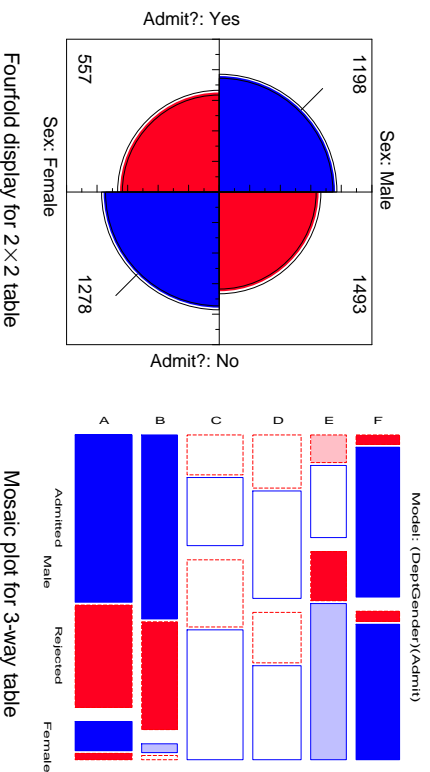
SUGI 28

7

Michael Friendly

Visual metaphors

- Quantitative data: **magnitude** \sim **position along an axis**
- Frequency data: **count** \sim **area**



SUGI 28

8

Michael Friendly

Categorical Data Analysis with Graphics

Exploratory methods

- Minimal assumptions (like non-parametric methods)
 - Show the data, not just *summaries*
 - Help detect *patterns*, *trends*, *anomalies*, suggest hypotheses
- ## Plots for model-based methods
- Residual plots - departures from model, omitted terms, ...
 - Effect plots - estimated probabilities of response or log odds
 - Diagnostic plots - influence, violation of assumptions
- ## Goals
- VCD and SSSG - Make these methods *available* and *accessible* in SAS
 - Practical power = Statistical power** \times **Probability of Use**
 - Today's goal: take-home knowledge
 - Tomorrow's goal: dynamic, interactive graphics for categorical data

SUGI 28

9

Michael Friendly

VCD Macros & SAS/IML programs

- Macros, datasets available at www.math.yorku.ca/SCS/vcd/

Discrete distributions

DISTPLOT Plots for discrete distributions
GOODFIT Goodness-of-fit for discrete distributions
ORDPLOT Ord plot for discrete distributions
POISSPLOT Poissonness plot
ROOTGRAM Hanging rootograms

Two-way and n-way tables

AGREE Observer agreement chart
CORRESP Plot PROC CORRESP results
FFOLD Fourfold displays for $2 \times 2 \times k$ tables (macro)
FOURFOLD Fourfold displays for $2 \times 2 \times k$ tables (SAS/IML)
SIEVE Sieve diagrams (SAS/IML)
MOSAIC Mosaic displays (macro)
MOSAICS SAS/IML modules for mosaic displays
MOSMAT Mosaic matrices (macro)
TABLE Construct a grouped frequency table, with recoding
TRIPLOT Trilinear plots for 72×3 tables

SUGI 28

10

Michael Friendly

Categorical Data Analysis with Graphics

Model-based methods

ADDVAR Added variable plots for logistic regression
CATPLOT Plot results from PROC CATMOD
HALFNORM Half-normal plots for generalized linear models
INFLGLIM Influence plots for generalized linear models
INFLLOGITS Influence plots for logistic regression
LOGODDS Plot empirical logits and probabilities for binary data
POWERLOG Power calculations for logistic regression
POWERKXC Power calculations for two-way frequency table
POWER2x2 Power calculations for a 2×2 table
ROBUST Robust fitting for linear models

Utility macros

DUMMY Create dummy variables
LAGS Calculate lagged frequencies for sequential analysis
PANELS Arrange multiple plots in a panelled display
SORT Sort a dataset by the value of a statistic or formatted value
 Utility
GRAPHICS Graphics utility macros: **BAR**, **EQUATE**, **GDISPLA**, **GENSYM**, **GSKIP**, **LABEL**, **POINTS**, **PSCALE**

VCD Archive (vcdprog.zip) available to purchasers at: support.sas.com/publishing/bnu/56571_sample.html

SUGI 28

11

Michael Friendly

Discrete distributions

- **Counts of occurrences:** accidents, words in text, blood cells with some characteristic.
- **Data:** Basic outcome value, k , $k = 0, 1, \dots$ and number of observations, n_k , with that value.

■ **Example:** distributions of key “marker” words: *from*, *may*, *whilst* . . . in *Federalist Papers* by James Madison, e.g., blocks of 200 words with *may*:

Occurrences (k)	0	1	2	3	4	5	6
Blocks (n_k)	156	63	29	8	4	1	1

k	0	1	2	3	4	5	6	7	8	9	10	11	12
n_k	3	24	104	286	670	1033	1343	1112	829	478	181	45	7

■ **Example:** Saxony families with 12 children having $k = 0, 1, \dots, 12$ sons.

SUGI 28

12

Michael Friendly

Discrete distributions

- **Questions:**
 - What process gave rise to the distribution?
 - Form of distribution: uniform, binomial, Poisson, negative binomial, geometric, etc.?
 - Estimate parameters
 - Visualize goodness of fit
- **For example:**
 - *Federalist Papers*: might expect a Poisson(λ) distribution.
 - *Families in Saxony*: might expect a Bin(n , p) distribution with $n = 12$. Perhaps $p = 0.5$ as well.

SUGI 28

13

Michael Friendly

Fitting and graphing discrete distributions

VCD methods to fit, visualize, and diagnose discrete distributions:

- **Fitting:** `GOODFIT` macro fits uniform, binomial, Poisson, negative binomial, geometric, logarithmic series distributions (or any specified multinomial)
- **Hanging rootograms:** Sensitive assess departure between Observed, Fitted counts (`ROOTGRAM` macro)
- **Ord plots:** Diagnose form of a discrete distribution (`ORDPLOT` macro)
- **Poissonness plots:** Robust fitting and diagnostic plots for Poisson (`POISPLIT` macro)
- **Robust distribution plots** (`DISTPLOT` macro)

SUGI 28

14

Michael Friendly

GOODFIT macro: Fitting discrete distributions

■ `GOODFIT` macro fits uniform, binomial, Poisson, negative binomial, geometric, logarithmic series distributions (or any specified multinomial)

■ E.g., Try fitting Poisson model

```

1 title "Instances of 'may' in Federalist papers";
2 data madison;
3   input count blocks;
4   label count='Number of Occurrences'
5     blocks='Blocks of Text';
6 datalines;
7   0      156
8   1      63
9   2      29
10  3      8
11  4      4
12  5      1
13  6      1
14 ;
15 %goodfit(data=madison, var=count, freq=blocks,
16   dist=poisson);

```

SUGI 28

15

Michael Friendly

Fitting discrete distributions

The **GOODFIT** macro gives a table of observed and fitted frequencies, Pearson χ^2 residuals (CHI) and likelihood-ratio deviance residuals (DEV).

Instances of 'may' in Federalist papers					
COUNT	BLOCKS	PHAT	EXP	CHI	DEV
0	156	0.51867	135.891	1.72499	6.56171
1	63	0.34050	89.211	-2.77509	-6.62056
2	29	0.11177	29.283	-0.05231	-0.75056
3	8	0.02446	6.408	0.62890	1.88423
4	4	0.00401	1.052	2.87493	3.26912
5	1	0.00053	0.138	2.31948	1.98992
6	1	0.00006	0.015	8.01267	2.89568
=====		=====		=====	
	262	0.99999	261.998		

SUGI 28

16

Michael Friendly

Fitting discrete distributions

In addition, it provides the overall goodness-of-fit tests:

```

Goodness-of-fit test for data set MADISON

Analysis variable:      COUNT Number of Occurrences
Distribution:           POISSON
Estimated Parameters:  lambda = 0.6565

Pearson chi-square     = 88.92304707
Prob > chi-square      = 0
Likelihood ratio G2    = 25.243121314
Prob > chi-square      = 0.0001250511
Degrees of freedom     = 5
  
```

The poisson model does not fit! Why?

SUGI 28

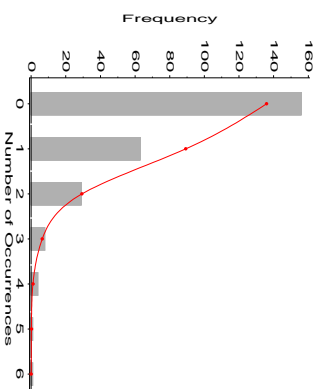
17

Michael Friendly

What's wrong with histograms?

- Discrete distributions often graphed as histograms, with a theoretical fitted distribution superimposed.

```
%goodfit(data=madison, var=count, freq=blocks,
dist=poisson);
```



Problems:

- largest frequencies dominate display
- must assess deviations vs. a curve

SUGI 28

18

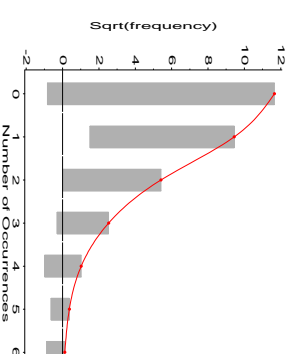
Michael Friendly

Hang & root them → Hanging rootograms

Tukey (1972, 1977):

- shift histogram bars to the fitted curve → judge deviations vs. horizontal line.
- plot $\sqrt{\text{freq}}$ → smaller frequencies are emphasized.

```
%goodfit(data=madison, var=count, freq=blocks,
dist=poisson, out=fit);
%rootgram(data=fit, var=count, obs=blocks);
```



SUGI 28

19

Michael Friendly

Ord plots: Diagnose form of discrete distribution

- How to tell which discrete distributions are likely candidates?
- Ord (1967): for each of Poisson, Binomial, Negative Binomial, and Logarithmic Series distributions,
 - plot of $k \cdot p_k / p_{k-1}$ against k is linear
 - signs of intercept and slope \rightarrow determine the form, give rough estimates of parameters

Slope (b)	Intercept (a)	Distribution (parameter)	Parameter estimate
0	+	Poisson (λ)	$\lambda \approx a$
-	+	Binomial (n, p)	$p \approx b / (b - 1)$
+	+	Neg. binomial (n, p)	$p \approx 1 - b$
+	-	Log. series (θ)	$\theta \approx b$ $\theta \approx -a$

- Fit line by WLS, using $\sqrt{n_k}$ as weights

SUGI 28

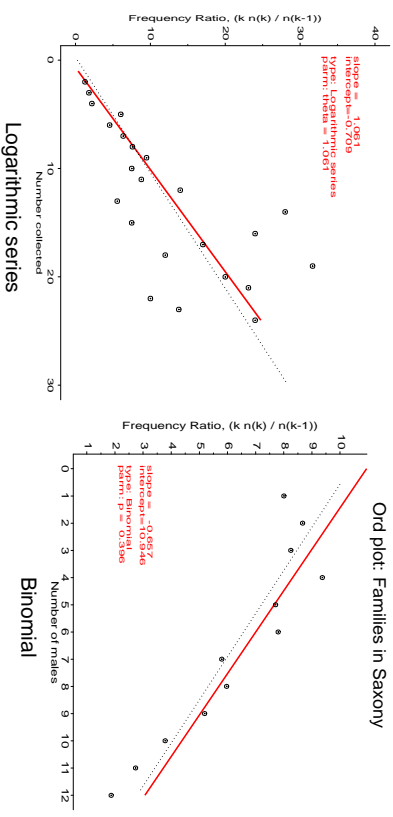
20

Michael Friendly

Ord plots

Other distributions:

Butterfly species collected in Malaysia



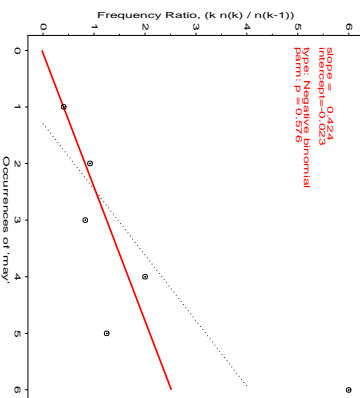
SUGI 28

22

Michael Friendly

Ord plots

- ORDPLOT** macro
 - `%ordplot (data=madison, count=Count, freq=blocks);`
 - Diagnoses distribution as NegBin
 - Estimates $\hat{p} = 0.576$



SUGI 28

21

Michael Friendly

Robust distribution plots: Poisson

- Ord plots lack robustness
 - one discrepant frequency, n_k affects points for both k and $k + 1$
 - Robust plots for Poisson distribution (Hoaglin and Tukey, 1985)
 - For Poisson, plot **count metameter** = $\phi(n_k) = \log_e(k! n_k / N^k)$ vs. k
 - Linear relation \Rightarrow Poisson, slope gives λ
 - CI for points, diagnostic (influence) plot
- POISPLLOT** macro

SUGI 28

23

Michael Friendly

POISPLOTT macro: example

```

1 title "Instances of 'may' in Federalist papers";
2 data madison;
3   input count blocks;
4   label count='Number of Occurrences'
5     blocks='Blocks of Text';
6   datalines;
7     0    156
8     1    63
9     2    29
10    3    8
11    4    4
12    5    1
13    6    1
14 ;
15 %poisplot(data=madison, count=count, freq=blocks);

```

Generalized robust distribution plots

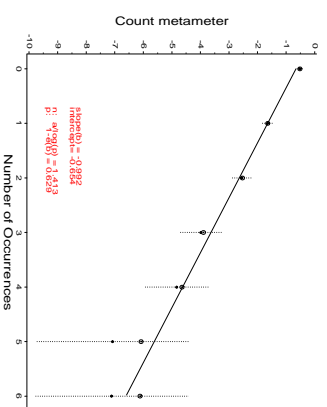
Other distributions: Analogous plots, for suitable count metamer, $\phi(n_k)$ vs. k .

- Linear relation \Rightarrow correct distribution, slope gives parameter estimates
- CI reflect variability of the individual counts, n_k
- **DISTPLOTT** macro

```

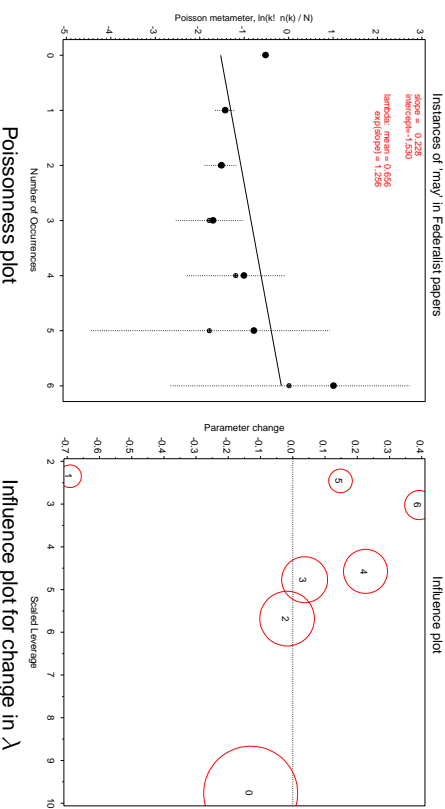
%distplot(data=madison, count=count, freq=blocks,
dist=negbin);

```



POISPLOTT macro: output

Curvilinear relation \rightarrow distribution is *not* Poisson



Contingency tables

- Two-way tables
 - 2×2 tables — Visualize odds ratio (**FFOLD** macro)
 - $2 \times 2 \times k$ tables — Homogeneity of association
 - $r \times 3$ tables — Trilinear plots (**TRIPLOT** macro)
 - $r \times c$ tables — Visualize association (**SIEVE** program)
 - $r \times c$ tables — Visualize association (**MOSAIC** macro)
 - Square $r \times r$ tables — Visualize agreement (**AGREE** program)
- r -way tables
 - Fit loglinear models, visualize lack-of-fit — (**MOSAIC** macro)
 - Test & visualize partial association — (**MOSAIC** macro)
 - Visualize pairwise association — (**MOSMAT** macro)
 - Visualize conditional association — (**MOSMAT** macro)
 - Visualize loglinear structure — (**MOSMAT** macro)
- Correspondence analysis and MCA — (**CORRESP** macro)

Methods for 2 × 2 tables

- Bickel et al. (1975): data on admissions to graduate departments at U. C. Berkeley in 1973.
- Aggregate data for the six largest departments:

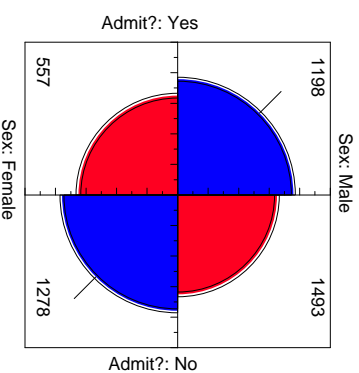
Table 3: Admissions to Berkeley graduate programs

	Admitted	Rejected	Total	% Admitted
Males	1198	1493	2691	44.52
Females	557	1278	1835	30.35
Total	1755	2771	4526	38.78

- Evidence for gender bias?
 - $G^2_{(1)} = 93.7, \chi^2_{(1)} = 92.2, p < 0.0001$
 - Odds ratio, $\theta = \frac{\text{Odds(Admit|Male)}}{\text{Odds(Admit|Female)}} = 1.84$
 - Males 84% more likely to be admitted.

Fourfold displays for 2 × 2 tables

- **Quarter circles:** radius $\sim \sqrt{n_{ij}} \Rightarrow \text{area} \sim \text{frequency}$
- **Independence:** Adjoining quadrants \approx align
- **Odds ratio:** ratio of areas of diagonally opposite cells
- **Confidence rings:** Visual test of $H_0: \theta = 1 \leftrightarrow$ adjoining rings overlap



- Confidence rings do not overlap: $\theta \neq 1$

Standard analysis: PROC FREQ

```
proc freq data=berkeley;
weight freq;
tables gender*admit / chisq;
```

Output:

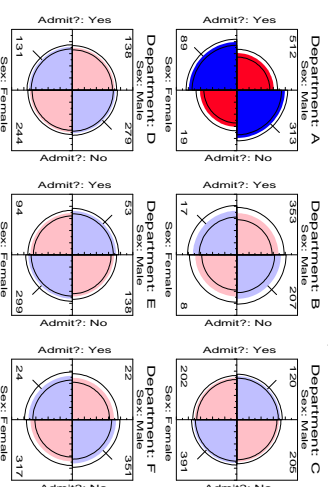
Statistics for Table of gender by admit

Statistic	DF	Value	Prob
Chi-Square	1	92.2053	< .0001
Likelihood Ratio Chi-Square	1	93.4494	< .0001
Continuity Adj. Chi-Square	1	91.6096	< .0001
Mantel-Haenszel Chi-Square	1	92.1849	< .0001
Phi Coefficient		0.1427	

How to visualize and interpret?

Fourfold displays for 2 × 2 × k tables

- Data in Table 3 had been pooled over departments
- Stratified analysis: one fourfold display for each department
- Each 2 × 2 table standardized to equate marginal frequencies
- Shading: highlight departments for which $H_{0i}: \theta_i \neq 1$



- Only one department (A) shows association: $\theta_A = 0.349 \rightarrow$ women $(0.349)^{-1} = 2.86$ times as likely as men to be admitted.

What happened here?

Simpson's paradox:

- Aggregate data are misleading because they falsely assume men and women apply *equally* in each field.
- But:
 - Large differences in admission rates across departments.
 - Men and women apply to these departments differentially.
 - Women applied in large numbers to departments with low admission rates.
- (This ignores possibility of structural bias against women: differential funding of fields to which women are more likely to apply.)
- Other graphical methods can show these effects.

SUGI 28

32

Michael Friendly

Two-way frequency tables: Sieve diagrams

- **count** ~ **area**
- When row/col variables are independent, $n_{ij} \sim n_{i+}n_{+j}$
- \Rightarrow each cell can be represented as a rectangle, with area = height \times width \sim frequency, n_{ij}

Expected frequencies: Hair Eye Color Data

	Green	Hazel	Blue	Brown	Black
Green	11.7	17.0	39.2	40.1	108
Hazel	30.9	44.9	103.9	106.3	286
Blue	7.7	11.2	25.8	26.4	71
Brown	13.7	20.0	46.1	47.2	127
Black	64	93	215	220	592
Hair Color					

SUGI 28

34

Michael Friendly

The FOURFOLD program and the FFOLD macro

- The **FOURFOLD** program is written in SAS/IML.
- The **FFOLD** macro provides a simpler interface.
- Printed output: (a) significance tests for individual odds ratios, (b) tests of homogeneity of association (here, over departments) and (c) conditional association (controlling for department).

Plot by department:

```
1 %include catdata(berkeley);
2
3 %ffold(data=berkeley,
4 var=Admit Gender,
5 by=Dept,
6 down=2, across=3,
7 htext=2);
```

berkf4f.sas

```
/* panel variables */
/* stratify by dept */
/* panel arrangement */
/* font size */
```

Aggregate data: first sum over departments, using the **TABLE** macro:

```
8 %table(data=berkeley, out=berk2,
9 var=Admit Gender,
10 weight=count,
11 order=data);
12 %ffold(data=berk2, var=Admit Gender);
```

/ omit dept */*
/ frequency variable */*

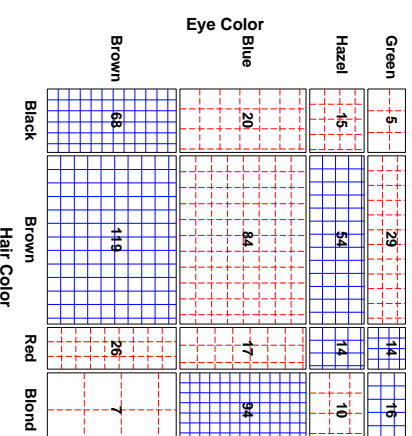
SUGI 28

33

Michael Friendly

Sieve diagrams

- Height/width ~ marginal frequencies, n_{i+}, n_{+j}
- Area ~ expected frequency, $\sim n_{i+}n_{+j}$
- Shading ~ observed frequency, n_{ij} ; color: sign($n_{ij} - \hat{n}_{ij}$).
- **Independence**: Shown when density of shading is uniform.



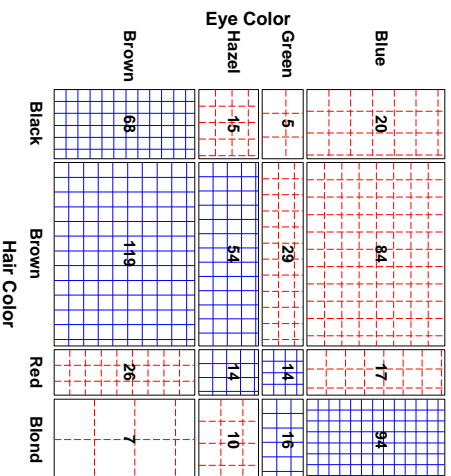
SUGI 28

35

Michael Friendly

Sieve diagrams

- Effect ordering: Reorder rows/cols to make the pattern coherent



Sieve diagrams: Example

Sieve2.sas

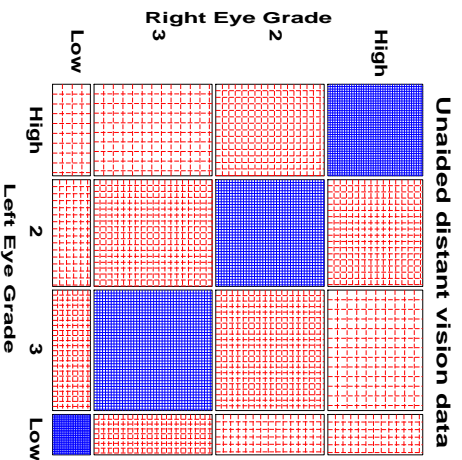
```

1 proc iml;
2   %include iml(sieve);
3   *--- frequency table;
4   tab = {1520 266 124 66,
5         234 1512 432 78,
6         117 362 1772 205,
7         36 82 179 492 };
8   *--- variable and level names;
9   vnames = {'Right Eye Grade' 'Left Eye Grade'};
10  Inames = {'High' '2' '3' 'Low',
11           'High' '2' '3' 'Low'};
12  title = {'Unaided distant vision data'};
13  *--- Global options;
14  font='hmpsl011';
15  run sieve(tab, vnames, Inames, title );
16  quit;

```

Sieve diagrams

- Vision classification data for 7477 women

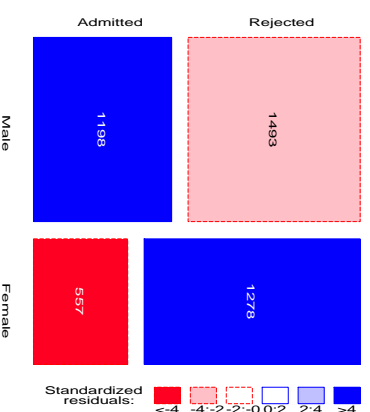


Mosaic displays and Log-linear Models

Hartigan and Kleiner (1981), Friendly (1994, 1999):

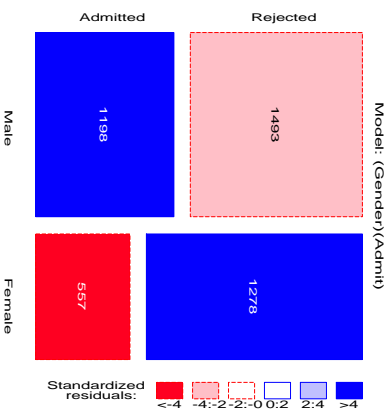
- Width ~ one set of marginals, n_{i+}
- Height ~ relative proportions of other variable, $p_{j|i} = n_{ij}/n_{i+}$
- area ~ frequency, $n_{ij} = n_{i+}p_{j|i}$

Model: (Gender)(Admit)



- **Shading:** Sign and magnitude of Pearson χ^2 residual, $d_{ij} = (n_{ij} - \hat{m}_{ij}) / \sqrt{\hat{m}_{ij}}$ (or L.R. G^2)
- Sign: — **negative in red**; + **positive in blue**
- Magnitude: intensity of shading; $|d_{ij}| > 0, 2, 4, \dots$

- **Independence:** Rows \approx align, or cells are empty!
- E.g., aggregate Berkeley data, independence model:



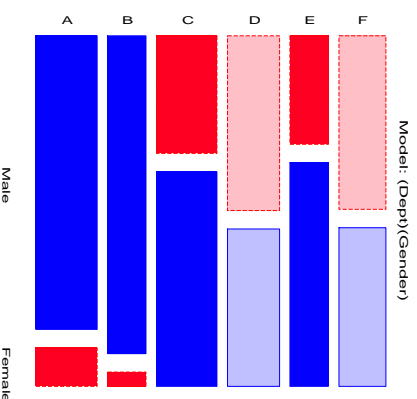
SUGI 28

40

Michael Friendly

Mosaic displays

- Departments \times Gender:
- Did departments differ in the total number of applicants?
 - Did men and women apply differentially to departments?



Model: (Dept)(Gender)

Model [Dept] [Gender]: $G^2_{(5)} = 1220.6$.

Note: Departments ordered A–F by overall rate of admission.

SUGI 28

41

Michael Friendly

Mosaic displays for multiway tables

- Generalizes to r -way tables: divide cells recursively
- Can fit any log-linear model (e.g., 3-way),

Table 4: Log-linear Models for Three-Way Tables

Model	Model symbol	Independence interpretation
Mutual independence	$[A][B][C]$	$A \perp B \perp C$
Joint independence	$[AB][C]$	$(AB) \perp C$
Conditional independence	$[AC][BC]$	$(A \perp B) C$
All two-way associations	$[AB][AC][BC]$	(none)
Saturated model	$[ABC]$	(none)

e.g., the model for conditional independence ($A \perp C | B$):

$$[AB][BC] \equiv \log m_{ijk} = \mu + \lambda_i^A + \lambda_j^B + \lambda_k^C + \lambda_{ij}^{AB} + \lambda_{jk}^{BC}$$

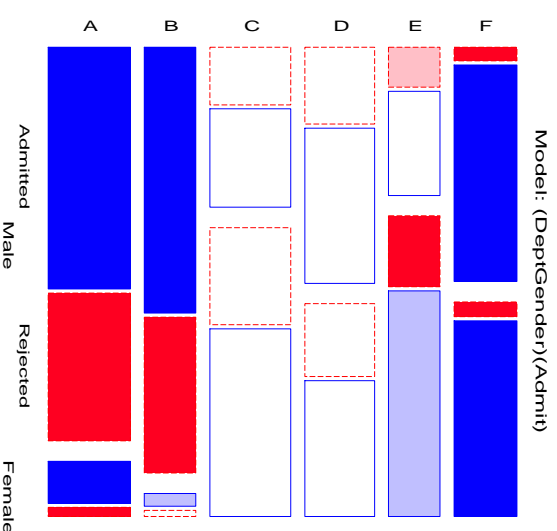
- Each mosaic shows:
- **DATA** (size of tiles)
- (some) **marginal** frequencies (spacing \rightarrow visual grouping)
- **RESIDUALS** (shading) — what associations have been omitted?

SUGI 28

42

Michael Friendly

- E.g., Joint independence (null model, Admit as response) [$G^2_{(11)} = 877.1$]:



SUGI 28

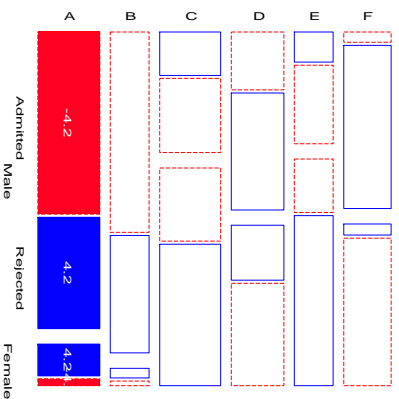
43

Michael Friendly

Mosaic displays for multiway tables

- Visual fitting:
 - Pattern of lack-of-fit (residuals) → "better" model—smaller residuals
 - "cleaning the mosaic" → "better" model—empty cells
 - best done interactively!

Model: (Dept|Gender)|(Dept|Admit)



- E.g., Add [Dept Admit] association → Conditional independence:
 - Fits poorly, overall ($G^2_{(6)} = 21.74$)
 - But, only in Department A!

SUGI 28

44

Michael Friendly

Categorical Data Analysis with Graphics

Software for Mosaic Displays

- Demonstration web applet:**
 - <http://www.math.yorku.ca/SCS/OnLine/mosaics/>
 - Runs the *current* version of mosaics via a cgi-script
 - Can run *sample data*, *upload* a data file, *enter* data in a form.
 - Choose model fitting and display options (not all supported).
- Documentation & software:**
 - <http://www.math.yorku.ca/SCS/mosaics.html>
- Examples:** Many in VCD and on web site
- SAS/IML modules:** *mosaics*, *sas* program
 - Enter frequency table directly in SAS/IML, or read from a SAS dataset.
 - Most flexible:
 - Select, collapse, reorder, re-label table levels using SAS/IML statements
 - Specify structural 0s, fit specialized models (e.g., quasi-independence)
 - Interface to models fit using PROC GENMOD

SUGI 28

45

Michael Friendly

Software for Mosaic Displays

- Macro interface:** *mosaic* macro, *table* macro, *mosmat* macro
- mosaic macro**
 - Easiest to use:
 - Direct input from a SAS dataset
 - No knowledge of SAS/IML required
 - Reorder table variables; collapse, reorder table levels with *table* macro
 - Convenient interface to *partial mosaics* (BY=)
- table macro**
 - Create frequency table from raw data
 - Collapse, reorder table categories
 - Re-code table categories using SAS formats, e.g., 1='Male' 2='Female'
- mosmat macro**
 - Mosaic matrices—analogy of scatterplot matrix (Friendly, 1999)

SUGI 28

46

Michael Friendly

Categorical Data Analysis with Graphics

Software for Mosaic Displays

- Other implementations:**
 - JMP and SAS/INSIGHT both provide rudimentary mosaic displays (two-way only, no interface with model-fitting engines (shame!).
 - The R-Project (<http://www.r-project.org>) now provides the *vcd* package, implementing most of the graphical methods from VCD.
 - Truly interactive mosaic displays have been implemented in:
 - Lisp-Stat (ViSta)—<http://forrest.psych.unc.edu/research/>
 - Java (Mondrian)—<http://http://www1.math.uni-augsburg.de/Wondrian/>

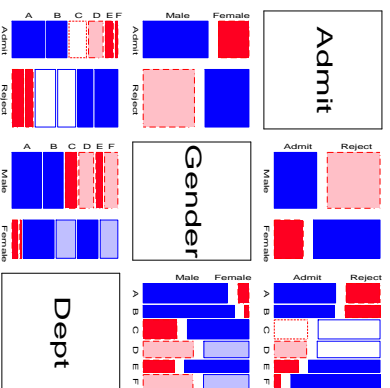
SUGI 28

47

Michael Friendly

mosmat macro: Mosaic matrices

```
1 %include catdata(berkeley);
2 %mosmat (data=berkeley,
3 vorder=Admit Gender Dept, sort=no);
```



Logit models

- **Fitting procedures**
 - PROC CATMOD
 - PROC LOGISTIC
 - PROC GENMOD / dist=poisson
 - SAS/INSIGHT (Fit Y X) Options → Distribution poisson
- **Visualization procedures**
 - CATPLOT macro - plot predicted, observed log odds from CATMOD
 - INFLGLIM macro - influence plots for generalized linear models
 - HALFNORM macro - half-normal plot of residuals for generalized linear models

Logit models

For a binary response, each loglinear model is equivalent to a logit model (logistic regression, with categorical predictors)

- Admit ⊥ Gender | Dept (conditional independence)

$$\log m_{ijk} = \mu + \chi_i^A + \chi_j^D + \chi_k^G + \chi_{ij}^{AD} + \chi_{jk}^{DG}$$

$$\leftrightarrow L_{ij} = \log(m_{ij1}/m_{ij2}) = \alpha + \beta_i^{\text{Dept}}$$

- Admit ⊥ Gender | Dept, except for Dept. A

$$\log m_{ijk} = \mu + \chi_i^A + \chi_j^D + \chi_k^G + \chi_{ij}^{AD} + \chi_{jk}^{DG} + \delta_{j=1} \chi_{ik}^{AG}$$

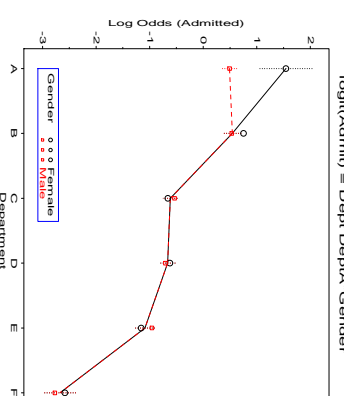
$$\leftrightarrow L_{ij} = \log(m_{ij1}/m_{ij2}) = \alpha + \beta_i^{\text{Dept}} + \delta_{j=1} \beta_{Gender}$$

where,

- $L_{ij} = \log(m_{ij1}/m_{ij2})$: log odds of admission for males as vs. females,
- β_i^{Dept} : effect on admissions of department,
- $\delta_{j=1} \beta_{Gender}$: effect of gender in Dept. A.

Plots for logit models

- Fit: PROC CATMOD; plot: CATPLOT macro
- Admit ⊥ Gender | Dept, except for Dept. A
proc catmod order=data data=berkeley;
...
response / out=predict;
model admit = dept dept|AG / m1;
%catplot (data=predict, xc=dept, class=gender,
type=FUNCTION, z=1.96, legend=Legend1);



Plots for logit models

- Fit: PROC CATMOD; plot: CATPLOT macro
- Admit ⊥ Gender | Dept, except for Dept. A

```

1 %include catdata(berkeley);
2 data berkeley;
3 set berkeley;
4 *-- Dummy variable for Gender in Dept A;
5 deptIAG = (gender='F') * (dept=1);
6 format dept dept.;
7
8 proc catmod order=data
9     data=berkeley;
10    weight freq;
11    population dept gender;
12    direct deptIAG;
13    response / out=predict;
14    model admit = dept deptIAG / ml;
15    run;
16    ...

```

SUGI 28

56

Michael Friendly

Plots for logit models

PROC CATMOD output:

Maximum Likelihood Analysis of Variance

Source	DF	Chi-Square	Pr > ChiSq
Intercept	1	291.22	<.0001
dept	5	571.45	<.0001
deptIAG	1	16.04	<.0001
Likelihood Ratio	5	2.68	0.7489

Analysis of Maximum Likelihood Estimates

Parameter	Estimate	Standard Error	Chi-Square	Pr > ChiSq
Intercept	-0.6685	0.0392	291.22	<.0001
dept A	1.1606	0.0705	271.21	<.0001
B	1.2113	0.0802	227.95	<.0001
C	0.0528	0.0687	0.59	0.4426
D	0.00358	0.0727	0.00	0.9607
E	-0.4210	0.0871	23.34	<.0001
deptIAG	1.0521	0.2627	16.04	<.0001

How to interpret?

SUGI 28

57

Michael Friendly

Plots for logit models

PROC CATMOD: observed and predicted logits:

```

17 proc print data=predict;
18     ... catberk6.sas ...
19 id dept gender;
20 var _obs_ _pred_ _sepred_;
21 format numeric_6.3 dept dept.;
22 where (_type_='FUNCTION');

```

dept	gender	_OBS_	_PRED_	_SEPPRED_
A	M	0.492	0.492	0.072
A	F	1.544	1.544	0.253
B	M	0.534	0.543	0.086
B	F	0.754	0.543	0.086
C	M	-0.536	-0.616	0.069
C	F	-0.660	-0.616	0.069
D	M	-0.704	-0.665	0.075
D	F	-0.622	-0.665	0.075
E	M	-0.957	-1.090	0.095
E	F	-1.157	-1.090	0.095
F	M	-2.770	-2.676	0.152
F	F	-2.581	-2.676	0.152

SUGI 28

58

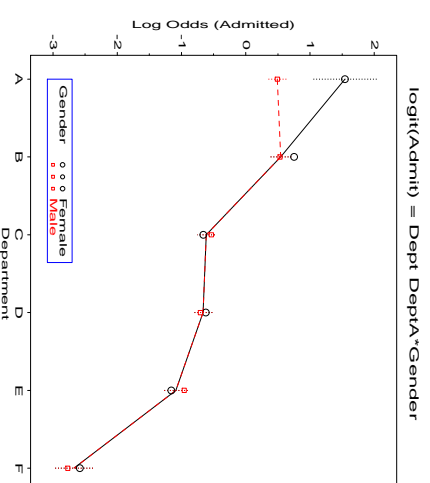
Michael Friendly

Plots for logit models

```

22 title 'Logit(Admit) = Dept DeptA*Gender';
23 %catplot(data=predict, x=dept, class=gender,
24 type=FUNCTION,
25 z=1.96);
26 *-- catberk6.sas
/* plot the log odds */
/* 95% error bars */

```



SUGI 28

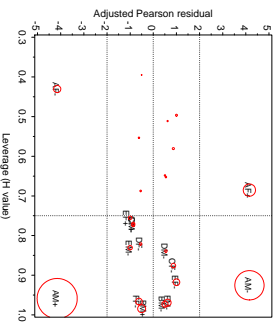
59

Michael Friendly

Diagnostic plots for Generalized Linear Models

INFLGLIM macro: Influence plots for generalized linear models (Williams, 1987)

- Fit: PROC GENMOD; calculates additional diagnostic measures (Hat value, Cook's D, etc.)
- Plot: measures of residual ($\hat{y} = \Delta \chi^2$, χ^2 residual) vs. leverage ($\hat{X} = \text{hat value}$), bubble size (area, radius) \sim Cook's D.
- which cells have undue impact on fitted model?



SUGI 28

60

Michael Friendly

Categorical Data Analysis with Graphics

INFLGLIM macro: Example

```

1 Berkeley data, model [AD][GD] ↔  $L_{ij} = \alpha + \beta_{\text{Dept}}$ 
2                               genberktl.sas
3 %include catdata(berkeley);
4 *-- make a cell ID variable, joining factors;
5 data berkeley;
6   set berkeley;
7   cell = trim(put(dept,dept.) ||
8     gender ||
9     trim(put(admit,yn.)));
10 %inflglm(data=berkeley,
11   class=dept gender admit,
12   resp=freq,
13   model=admit|dept gender|dept,
14   dist=poisson,
15   id=cell,
16   gx=hat, gy=streschi);

```

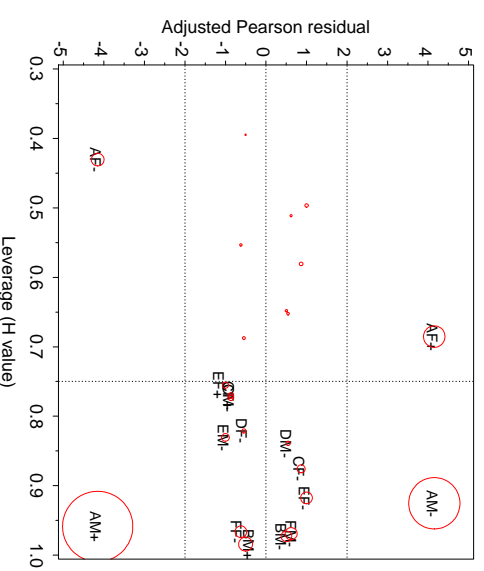
SUGI 28

61

Michael Friendly

INFLGLIM macro: Example

- All cells which do not fit ($|r_i^*| > 2$) are for department A.
- Males applying to dept A have large leverage \Rightarrow large influence (Cook's D)



SUGI 28

62

Michael Friendly

Categorical Data Analysis with Graphics

Diagnostic plots for Generalized Linear Models

HALFNORM macro: Half-normal plot of residuals (Atkinson, 1981)

- Plot ordered absolute residuals, $|r_i^{(a)}|$ vs. expected normal values, $|z_i^{(a)}|$
- Standard normal confidence envelope not suitable for GLMs
- Simulate reference 'line' and envelope with simulated confidence intervals

```

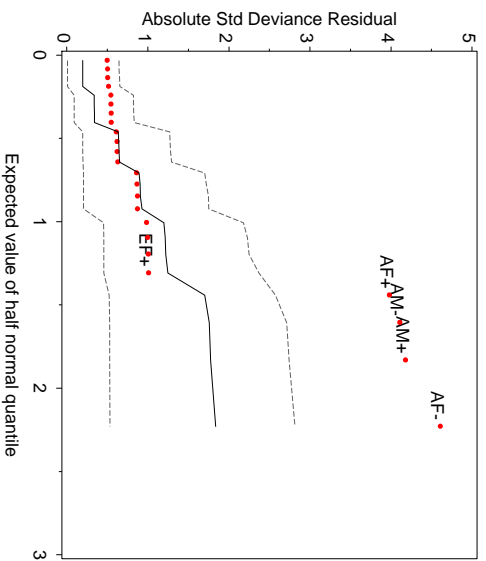
1 %halfnorm(data=berkeley,
2   class=dept gender admit,
3   resp=freq,
4   model=dept|gender dept|admit,
5   dist=poisson, id=cell);
6   ... genberktl.sas

```

SUGI 28

63

Michael Friendly



- Points with largest |residual| labeled
- The model fits well, except in department A.

SUGI 28

64

Michael Friendly

Logistic regression models

- **Response variable:**
 - Binary response: success/failure, vote: yes/no
 - Ordinal response: none, some, severe depression
 - Polytomous response: vote Liberal, Tory, Alliance, NDP
- **Explanatory variables:**
 - Quantitative regressors: age, dose
 - Transformed regressors: $\sqrt{\text{age}}$, $\log(\text{dose})$
 - Polynomial regressors: age^2 , age^3 , ...
 - Categorical predictors: treatment, sex
 - Interaction regressors: treatment \times age, sex \times age

SUGI 28

65

Michael Friendly

Logistic regression models: Binary response

- For a binary response, $Y \in (0, 1)$, let x be a vector of p regressors, and π_i be the probability, $P_i(Y = 1 | x)$.
- The logistic regression model is a linear model for the *log odds*, or *logit* that $Y = 1$, given the values in x ,

$$\begin{aligned} \text{logit}(\pi_i) &\equiv \log\left(\frac{\pi_i}{1 - \pi_i}\right) = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} \\ &= \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} \end{aligned}$$

- An equivalent (non-linear) form of the model may be specified for the probability, π_i , itself,

$$\pi_i = \{1 + \exp(-[\alpha + \mathbf{x}_i^T \boldsymbol{\beta}])\}^{-1}$$

- The logistic model is a *linear model* for the log odds, but also a *multiplicative model* for the odds of “success,”

$$\frac{\pi_i}{1 - \pi_i} = \exp(\alpha + \mathbf{x}_i^T \boldsymbol{\beta}) = \exp(\alpha) \exp(\mathbf{x}_i^T \boldsymbol{\beta})$$

so, increasing x_{ij} by 1 increases $\text{logit}(\pi_i)$ by β_j , and multiplies the odds by e^{β_j} .

SUGI 28

66

Michael Friendly

Logistic regression models: Binary response

- **Fitting:** PROC LOGISTIC (or ROBUST macro—M-estimation)
 - Data:
 - Frequency form (from PROC FREQ)—when all predictors are discrete
 - Case form—when any predictors are quantitative
 - Models:
 - CLASS statement (V7+)—no need for dummy variables
 - discrete predictors
 - can specify *order* and *parameterization* (effect, polynomial, reference cell)
- MODEL statement—allows GLM syntax, e.g.,
model Better = Sex — Treat — Age @2;

SUGI 28

67

Michael Friendly

Logistic regression models: Binary response

Visualization:

- Goal: see and *understand* the data and fitted model
- **LOGODDS** macro: Plot observed responses, fitted and smoothed probabilities
- Model plots:
 - OUTPUT statement →
 - fitted $\hat{\pi}_i$, lower/upper $(1 - \alpha)$ CI, and/or
 - fitted logit $(\alpha + x_i^T \beta) \pm z_{1-\alpha/2} \text{se}(\text{logit})$
 - Plot with standard procedures (PROC GCHART, GPLOT)
- Utility macros (**BARS**, **LABEL**, **POINTS**, **PSCALE**, etc.) for custom displays
- Effect plots — plot hierarchical subset of effects, averaging over those not included.
- **INFLOGITS** macro: Influence plots for logistic regression models
- **ADDVAR** macro: Added variable plots for new predictors or transformations of old

SUGI 28

68

Michael Friendly

Example: Arthritis treatment data

- Predictors: Sex, Treatment (treated, placebo), Age
- Response: improvement (none, some, marked)
- Consider first as binary response: None vs. (Some or Marked)=‘Better’

```

1 Data in case form: arthritis.sas
2 data arthritis;
3   length treat $7. sex $6. ;
4   input id treat $ sex $ age improve @@ ;
5   better = (improve > 0); *-- Make binary response;
6 datalines ;
7 57 Treated Male 27 1 9 Placebo Male 37 0
8 46 Treated Male 29 0 14 Placebo Male 44 0
9 77 Treated Male 30 0 73 Placebo Male 50 0
10 ..... (observations omitted)
11 56 Treated Female 69 1 42 Placebo Female 66 0
12 12 Treated Female 15 1 Placebo Female 66 1
13 71 Treated Female 71 1 Placebo Female 68 1
14 1 Placebo Female 74 2
15 ;
  
```

SUGI 28

69

Michael Friendly

LOGODDS macro: Empirical logit plots

- **Linearity**: Is a linear relation realistic?
 - **Smoothing**: Discrete data often requires smoothing to see!
- The **LOGODDS** macro:

- Show the data: Plot (0/1) responses (stacked or jittered)
- Divide X into groups (e.g., deciles), empirical logit, $\text{log} \left(\frac{y_i+1/2}{n_i-y_i+1/2} \right)$, for each
- Linear logistic regression, plus smoothed curve (**LOWESS** macro)

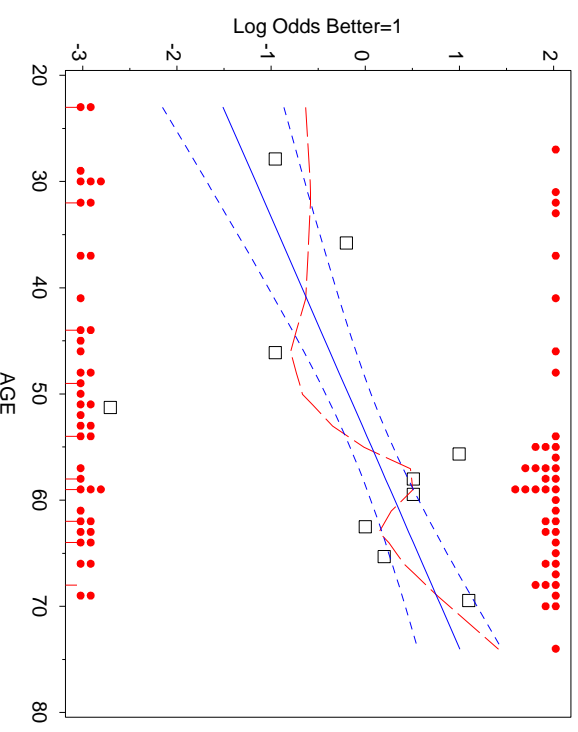
```

1 %include catdata(arthritis);
2 %logods (data=arthritis,
3   x=age, y=Better, /* vars to plot
4   smooth=0.5, /* LOWESS smoothing parameter */
5   plot=logit); /* plot on logit scale */
  
```

SUGI 28

70

Michael Friendly



SUGI 28

71

Michael Friendly

PROC LOGISTIC: Model fitting and plotting

- Specify ordering of response levels (order= or descending options)
- Specify parameterizations for CLASS variables
- OUTPUT statement to get fitted logits and probabilities

```

1 logistic.sas ...
2 proc logistic data=arthritis descending;
3   class sex (ref=last) treat (ref=first) / param=ref;
4   model better = sex treat age;
5   output out=results
6     p=prob l=lower u=upper
7     xbeta=logit stdxbeta=selogit / alpha=.33;

```

The output includes:

Effect	DF	Chi-Square	Wald		Pr > ChiSq
			Chi-Square	Pr > ChiSq	
sex	1	6.2576	6.2576	0.0124	0.0124
treat	1	10.7596	10.7596	0.0010	0.0010
age	1	5.5655	5.5655	0.0183	0.0183

SUGI 28

72

Michael Friendly

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error		Chi-Square	Pr > ChiSq
			Estimate	Error		
Intercept	1	-4.5033	1.3074	11.8649	0.0006	
sex	1	1.4878	0.5948	6.2576	0.0124	
treat	1	1.7598	0.5365	10.7596	0.0010	
age	1	0.0487	0.0207	5.5655	0.0183	

Odds Ratio Estimates

Effect	Point Estimate	95% Wald Confidence Limits	
		Estimate	Limits
sex	Female vs Male	4.427	1.380 14.204
treat	Treated vs Placebo	5.811	2.031 16.632
age		1.050	1.008 1.093

Parameter estimates (reference cell coding):

- $\beta_1 = 1.49 \Rightarrow$ Females $e^{1.49} = 4.43$ times more likely to be better than Males
- $\beta_2 = 1.76 \Rightarrow$ Treated $e^{1.76} = 5.81$ times more likely to be better than Placebo
- $\beta_3 = 0.0487 \Rightarrow$ odds ratio = 1.05 \Rightarrow odds of improvement increase 5% each year. Over 10 years, odds of improvement multiplied by $e^{10 \times 0.0487} = 1.63$, a 63% increase.

SUGI 28

73

Michael Friendly

PROC LOGISTIC: Model plots

- Plots of fitted values from the dataset specified on the OUTPUT statement
- Plot either predicted probabilities or logits

The first few observations from the results dataset:

id	sex	treat	age	better	prob	lower	upper	Logit	selogit
57	Male	Treated	27	1	0.194	0.103	0.334	-1.427	0.758
9	Male	Placebo	37	0	0.063	0.032	0.120	-2.700	0.725
46	Male	Treated	29	0	0.209	0.115	0.350	-1.330	0.728
14	Male	Placebo	44	0	0.086	0.047	0.152	-2.358	0.658
77	Male	Treated	30	0	0.217	0.122	0.357	-1.281	0.713
73	Male	Placebo	50	0	0.112	0.065	0.188	-2.066	0.622
...									

SUGI 28

74

Michael Friendly

PROC LOGISTIC: Model plots

Basic plots:

- Plot either logit or probability vs. one predictor (continuous or most levels)
- Separate curves for one factor
- Separate panels for all others (BY statement)

```

proc gplot data=results;
  plot (logit prob) * age = treat;
  by sex;
  symbol1 v=circle i=join l=3 c=black;
  symbol2 v=dot i=join l=1 c=red;

```

SUGI 28

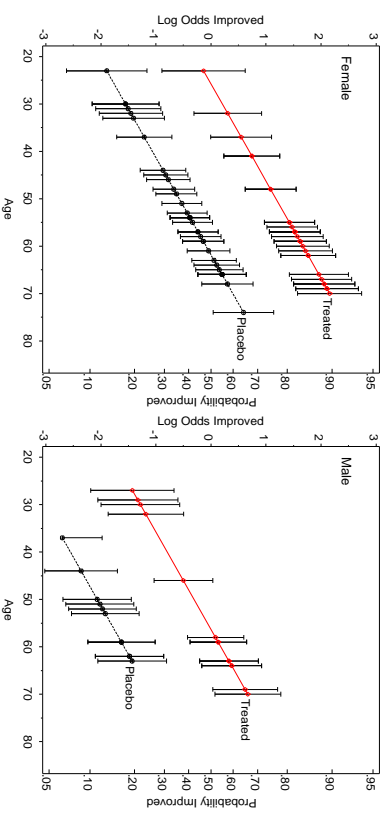
75

Michael Friendly

PROC LOGISTIC: Model plots

Enhanced plots:

- Plot on logit scale, with probability scale at right (PSCALE macro)
- Show 67% error bars $\approx \pm 1$ se (BARS macro)
- Custom legend and panel labels (LABEL macro)



SUGI 28

76

Michael Friendly

Categorical Data Analysis with Graphics

PROC LOGISTIC: Model plots

Enhanced plots:

... glogistic.sas ...

```

9 *** Error bars, on logit scale;
10 %bars(data=results, var=logit,
11 class=age, cvar=treat, by=age,
12 barlen=selogit, out=bars);
13
14 *** Custom legends and panel labels;
15 %label(data=results, y=logit, x=age, koff=1, cvar=treat,
16 by=sex, subset=last.treat, out=label1, pos=6, text=treat);
17 %label(data=results, y=2.5, x=20, size=2,
18 by=sex, subset=first.sex, out=label2, pos=6, text=sex);
19
20 *** Probability scales at right;
21 %pscale(out=pscale,
22 byvar=sex, byval=%str('Female', 'Male'));
23
24 data bars;
25 set label1 label2 bars pscale;
26 proc sort;
27 by sex;
28

```

SUGI 28

77

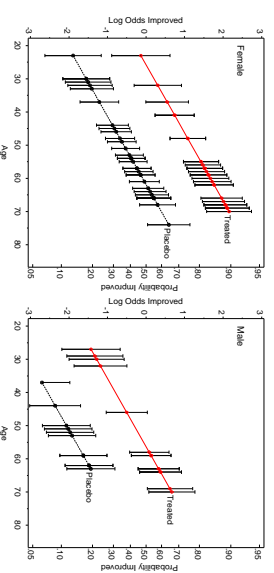
Michael Friendly

... glogistic.sas

```

30 title ' ';
31 h=1.8 a=-90 'Probability Improved'
32 h=2.5 a=-90 ' ';
33 goptions hby=0;
34 proc gplot data=results;
35 plot logit * age = treat /
36 vaxis=axis1 haxis=axis2 hm=1 vm=1
37 nolegend anno=bars frame;
38 by sex;
39 axis1 label=(a=90 'Log Odds Improved');
40 order=(3 to 3);
41 axis2 order=(20 to 80 by 10) offset=(2,6);
42 symbol1 v=+ i=join l=3 c=black;
43 symbol2 v=+ i=join l=1 c=red;
44 label age= 'Age' ;
45 run;

```



SUGI 28

78

Michael Friendly

Categorical Data Analysis with Graphics

Models with interactions

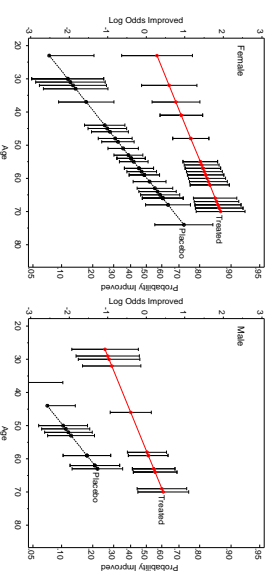
Plotting fitted values

- Only need to change the MODEL statement
- Output dataset automatically incorporates all model terms
- Plotting steps remain exactly the same

```

1 proc logistic data=artbrit descending;
2 class sex (ref=last) treat (ref=first) / param=ref;
3 model better = sex treat | age @2;
4 output out=results p=prob l=lower u=upper
5 xbeta=logit stdxbeta=selogit / alpha=.33;

```



SUGI 28

79

Michael Friendly

Effect plots for generalized linear models

- Fox (1987) — For complex models, often wish to plot a specific main effect or interaction (including lower-order relatives)
 - Fit full model to data with linear predictor (e.g., logit) $\eta = \mathbf{X}\beta$ and link function $g(\mu) = \eta \rightarrow$ estimate \mathbf{b} of β and covariance matrix $V(\mathbf{b})$ of \mathbf{b} .
 - Vary each predictor in the term over its' range
 - Fix other predictors at "typical" values (mean, median, proportion in the data) \rightarrow "effect model matrix," \mathbf{X}^*
 - Calculate fitted effect values, $\hat{\eta}^* = \mathbf{X}^* \mathbf{b}$.
 - Standard errors are square roots of $\text{diag}(\mathbf{X}^* V(\mathbf{b}) \mathbf{X}^{*T})$
 - Plot $\hat{\eta}^*$, or values transformed back to scale of response, $g^{-1}(\hat{\eta}^*)$.
- Note: This provides a general means to visualize interactions in all linear and generalized linear models.

SUGI 28

80

Michael Friendly

Effect plots in SAS

- Create a grid of values for predictors in the effect (EXPGRID macro)
- Fix other predictors at "typical" values (mean, median, proportion in the data)
- Concatenate grid with data
- Fit model \rightarrow output data set \rightarrow fitted values in the grid
- Standard errors automatically calculated
- Plot fitted values in the grid
- (Not yet a macro)

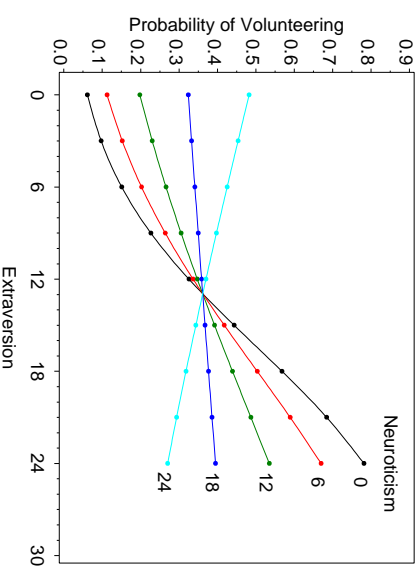
SUGI 28

81

Michael Friendly

Effect plots: Example

- Cowles and Davis (1987) — Volunteering for a psychology experiment
- Predictors: Sex, Neuroticism, Extraversion
- \rightarrow strong interaction, Neuroticism \times Extraversion



SUGI 28

82

Michael Friendly

Effect plots: Example

```

1 %include catdata(cowles);
2 %expgrid(Sex=0.5, /* Fix Sex at mid value */
3 Extraver=0 to 24 by 3, /* range of predictors */
4 Neurot=0 to 24 by 6, /* grid select value */
5 _in_=-1);
6
7 *** concatenate grid values with data;
8 data cowles;
9 set cowles _grid_;
10
11 *** fit model, output fitted values;
12 proc logistic data=cowles outest=parm covout;
13 model Volunter = Sex Extraver | Neurot / covb;
14 output out=predicted xbeta=logit stdxbeta=selogit
15 p=prob u=upper l=lower / alpha=.35;
16
17 *** select grid, replace labels;
18 data effect;
19 set predicted (where =(_in_=-1));
20 label logit='Log odds of Volunteering'
21 prob = 'Probability of Volunteering';

```

SUGI 28

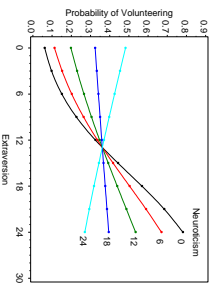
83

Michael Friendly

Effect plots: Example

```

23 *--- Custom legend;
24 %label (data=effect,
25 x=Extraver, y=prob,
26 subset=Extraver=24,
27 text=put(Neurot,3.),
28 pos=6, xoff=.2, out=labels);
29
30 *--- Plot step;
31 proc gplot data=effect;
32 plot prob * Extraver = Neurot /
33     vaxis=axis1 haxis=axis2 vm=1
34     anno=labels nolegend;
35     symbol v=dot i=spline r=5;
36     axis1 label=(a=90 r=0) order=(0 to .9 by .1);
37     axis2 order=(0 to 30 by 6) offset=(3,1);
38 run; quit;
    
```



SUGI 28

84

Michael Friendly

Influence measures and diagnostic plots

- **Leverage:** Potential impact of an individual case ~ distance from the centroid in space of predictors
- **Residuals:** Which observations are poorly fitted?
- **Influence:** Actual impact of an individual case ~ leverage × residual
- **C. CBAR** – analogs of Cook's D in OLS ~ standardized change in regression coefficients when i -th case is deleted.
- **DIFCHISQ, DIFDEV** – $\Delta\chi^2$ when i -th case is deleted.

SUGI 28

85

Michael Friendly

Influence measures and diagnostic plots

```

1 PROC LOGISTIC provides printed output with the influence and iplots options
2 model better = sex treat age / influence iplots;
    
```

Case Number	Value	The LOGISTIC Procedure Deviance Residual				Hat Matrix Diagonal											
		(1 unit = 0.26)	-8	-4	0	2	4	6	8	(1 unit = 0.01)	0	2	4	6	8	12	16
1	1.812	*								0.089							*
2	0.360		*							0.031							*
3	0.685			*						0.087							*
4	0.425				*					0.034							*
5	0.700					*				0.086							*
6	0.488						*			0.038							*
7	1.703	*								0.084							*
8	0.499				*					0.039							*
9	1.396					*				0.066							*
10	0.511						*			0.040							*
11	1.142				*					0.064							*
12	0.523					*				0.041							*
13	1.234						*			0.065							*
14	0.599				*					0.051							*
15	1.121		*							0.065							*
16	0.599			*						0.051							*

SUGI 28

86

Michael Friendly

17	1.319				*					0.069							*
18	0.640				*					0.058							*
19	1.319				*					0.069							*
20	0.640				*					0.058							*
21	1.340				*					0.070							*
22	1.814	*								0.061							*
23	1.022		*							0.070							*
24	0.529			*						0.060							*
25	1.449				*					0.078							*
26	0.619			*						0.053							*
27	0.909		*							0.080							*

- Problems:
- Way too much output
 - Doesn't highlight unusual cases well
 - Index plots don't consider combinations of measures

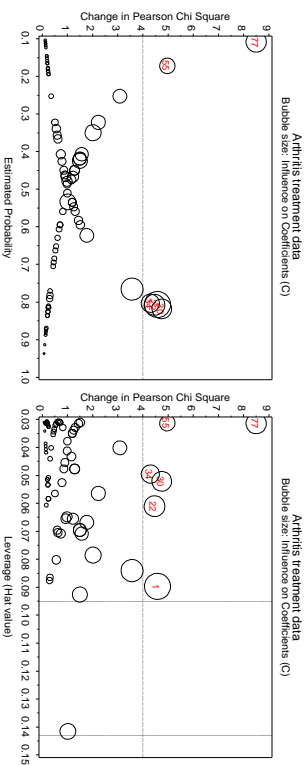
SUGI 28

87

Michael Friendly

INFLOGIS macro

- Specialized version of **INFLOGITM** macro for logistic regression
- Plots a measure of change in χ^2 (DIFCHISQ or DIFDEV) vs. predicted probability or leverage.
- Bubble symbols show actual influence (C or CBAR)
- Shows standard cutoffs for "large" values
- Labels outlying cases



SUGI 28

88

Michael Friendly

Categorical Data Analysis with Graphics

INFLOGIS macro: Example

```

1 %include data(arthrit);
2 %inflogis(data=arthrit,
3 class=sex treat,
4 y=better,
5 x=sex treat age,
6 id=case,
7 gy=DIFCHISQ,
8 gx=PRD HAT);

```

/ CLASS variables */*
/ response */*
/ predictors */*
/ case ID */*
/ graph ordinate */*
/ graph abscissas */*

Printed output lists cases with "large" leverage, residual or influence:

case	better	sex	treat	age	pred	hat	difchisq	difdev	c
1	1	Male	Treated	27	.806	.09	4.578	3.695	0.451
22	1	Male	Placebo	63	.807	.06	4.460	3.565	0.290
30	1	Female	Placebo	31	.818	.05	4.749	3.657	0.261
34	1	Female	Placebo	33	.803	.05	4.296	3.464	0.224
55	0	Female	Treated	58	.172	.03	4.970	3.676	0.160
77	0	Female	Treated	69	.108	.03	8.498	4.712	0.276

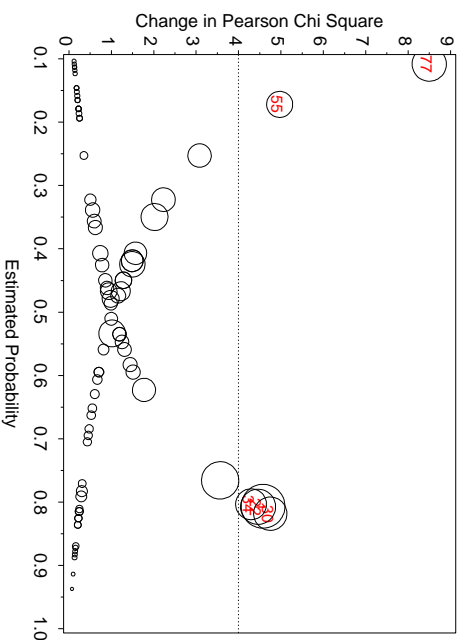
SUGI 28

89

Michael Friendly

INFLOGIS macro: Example

Arthritis treatment data
Bubble size: Influence on Coefficients (C)



SUGI 28

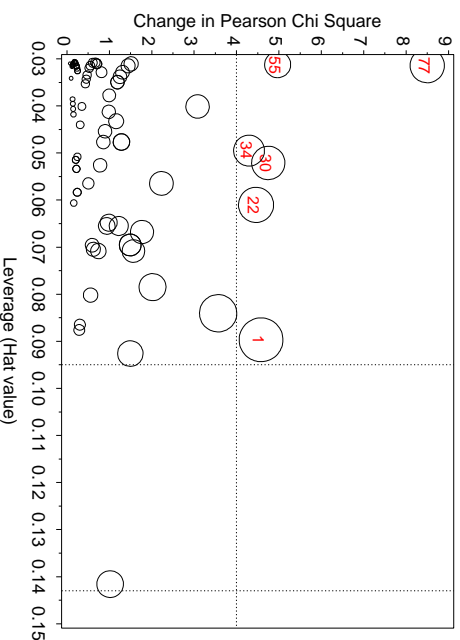
90

Michael Friendly

Categorical Data Analysis with Graphics

INFLOGIS macro: Example

Arthritis treatment data
Bubble size: Influence on Coefficients (C)



SUGI 28

91

Michael Friendly

Conclusions

- **Summarization & exposure**
 - Effective data analysis requires *summarization*—hypothesis tests, model fits (& comparisons!), parameter estimates (& precision)
 - Also requires *exposure*—displays to help the viewer see (& understand!) patterns, trends, and anomalies.
- **Graphical methods for categorical data**
 - Many new methods developed over the last 10–15 years
 - Some novel, others extend familiar methods for quantitative data
 - Described and illustrated in VCD
- **Theory into practice**
 - To be useful, statistical methods must be:
 - available—implemented in standard software
 - accessible—easy to use (or at least easier)
 - VCD provides ~40 general macros and SAS/IML programs

SUGI 28

Categorical Data Analysis with Graphics

92

Michael Friendly

References

- Atkinson, A. C. Two graphical displays for outlying and influential observations in regression. *Biometrika*, 68:13–20, 1981.
- Bickel, P. J., Hammel, J. W., and O'Connell, J. W. Sex bias in graduate admissions: Data from Berkeley. *Science*, 187:396–403, 1975.
- Cowles, M. and Davis, C. The subject matter of psychology: Volunteers. *British Journal of Social Psychology*, 26:97–102, 1987.
- Fox, J. Effect displays for generalized linear models. In Clogg, C. C., editor, *Sociological Methodology*, 1987, pp. 347–361. Jossey-Bass, San Francisco, 1987.
- Friendly, M. Mosaic displays for multi-way contingency tables. *Journal of the American Statistical Association*, 89:190–200, 1994.
- Friendly, M. Extending mosaic displays: Marginal, conditional, and partial views of categorical data. *Journal of Computational and Graphical Statistics*, 8(3):373–395, 1999.
- Friendly, M. Multidimensional arrays in SAS/IML. In *Proceedings of the SAS Users Group International Conference*, volume 25, pp. 1420–1427. SAS Institute, 2000.
- Friendly, M. Corgrams: Exploratory displays for correlation matrices. *The American Statistician*, 56(4):316–324, 2002.
- Friendly, M. and Kwan, E. Effect ordering for data displays. *Computational Statistics and Data Analysis*, 37, 2002. In press.
- Hartigan, J. A. and Kleinier, B. Mosaics for contingency tables. In Eddy, W. F., editor, *Computer Science and Statistics: Proceedings of the 13th Symposium on the Interface*, pp. 268–273. Springer-Verlag, New York, NY, 1981.

SUGI 28

93

Michael Friendly

- Hoaglin, D. C. and Tukey, J. W. Checking the shape of discrete distributions. In Hoaglin, D. C., Mosteller, F., and Tukey, J. W., editors, *Exploring Data Tables, Trends and Shapes*, chapter 9. John Wiley and Sons, New York, 1985.
- Ord, J. K. Graphical methods for a class of discrete distributions. *Journal of the Royal Statistical Society, Series A*, 130:232–238, 1967.
- Tufts, E. R. *The Visual Display of Quantitative Information*. Graphics Press, Cheshire, CT, 1983.
- Tukey, J. W. Some graphic and semigraphic displays. In Bancroft, T. A., editor, *Statistical Papers in Honor of George W. Snedecor*. Iowa State University Press, Ames, IA, 1972.
- Tukey, J. W. *Exploratory Data Analysis*. Addison Wesley, Reading, MA, 1977.
- Williams, D. A. Generalized linear model diagnostics using the deviance and single case deletions. *Applied Statistics*, 36:181–191, 1987.

SUGI 28

94

Michael Friendly