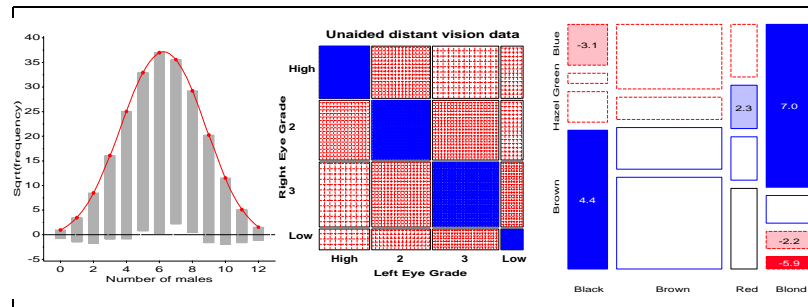


## Categorical Data Analysis with Graphics



Michael Friendly

York University, <friendly@yorku.ca>



SCS Short Course

Feb-Mar., 2007



## 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; SPSS: Crosstabs
  - 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  $\Pr\{\text{response}\}$
- SPSS: Hiloglinear, Loglinear, Generalized linear models

## Outline

- Overview: Categorical Data *and* Graphics
- Methods for discrete distributions – testing goodness of fit
  - Hanging rootograms
  - Robust distribution plots
- Methods for two-way frequency tables – understanding association
  - Fourfold displays
  - Sieve diagrams
- Mosaic displays and loglinear models for  $n$ -way tables
  - Mosaic displays
  - Mosaic matrices
  - Correspondence analysis and MCA
- Logistic and logit regression
  - Logit plots, effect plots
  - Diagnostic plots

Color version of these slides:

<http://www.math.yorku.ca/SCS/Courses/grcat/>

## 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.*

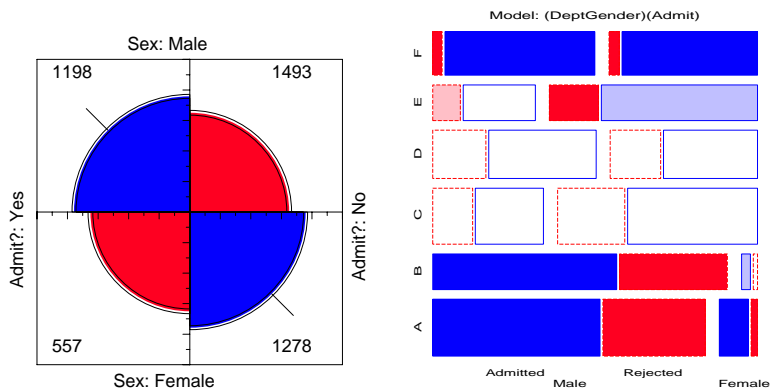
Farquhar & Farquhar, 1891

### ■ 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?

■ Visual metaphors

- Quantitative data: **magnitude** ~ **position along an axis**
- Frequency data: **count** ~ **area** (Friendly, 1995)



Fourfold display for 2x2 table

Mosaic plot for 3-way table

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

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

Eye color	Hair color			
	Blond	Black	Brown	Red
Blue	94	20	17	84
Brown	7	68	26	119
Green	10	15	14	54
Hazel	16	5	14	29

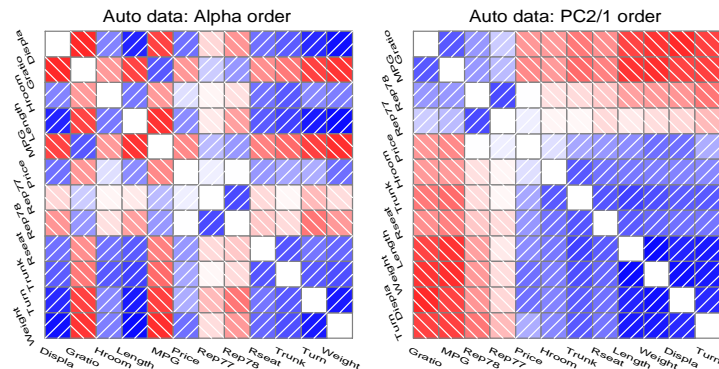
Table 2: Hair color - Eye color data: Effect ordered

Eye color	Hair color			
	Black	Brown	Red	Blond
Brown	68	119	26	7
Hazel	15	54	14	10
Green	5	29	14	16
Blue	20	84	17	94

Model:	Independence: [Hair][Eye] $\chi^2(9) = 138.29$							
Color coding:	<-4	<-2	<-1	0	>1	>2	>4	
n in each cell:	n < expected				n > expected			

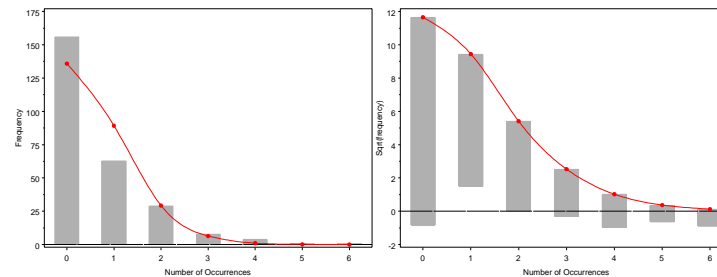
■ Principles of Graphical Displays

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



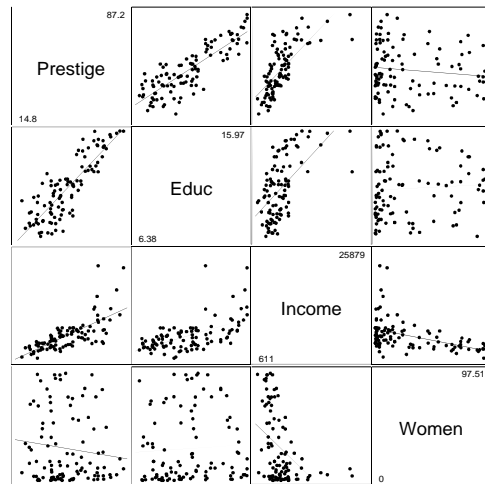
"Corrgrams: Exploratory displays for correlation matrices" (Friendly, 2002)

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



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

- e.g., scatterplot matrix for quantitative data: all pairwise scatterplots



■ **Exploratory methods**

- Minimal assumptions (like non-parametric methods)
- Show the *data*, not just *summaries*
- Help detect *patterns, trends, anomalies*, suggest hypotheses

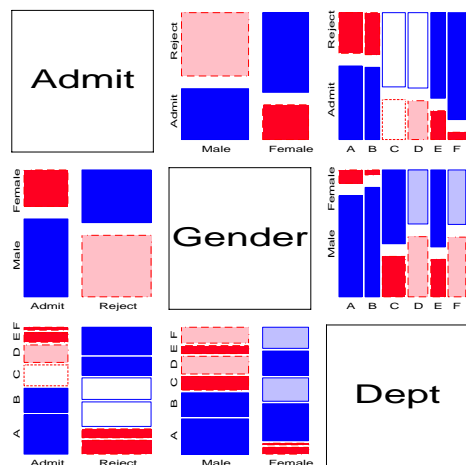
■ **Plots for model-based methods**

- 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 × Probability of Use**
- Today's goal: take-home knowledge
- Tomorrow's goal: dynamic, interactive graphics for categorical data

- e.g., mosaic matrix for quantitative data: all pairwise mosaic plots



**VCD Macros & SAS/IML programs**

- Macros, datasets available at [www.math.yorku.ca/SCS/vcd/](http://www.math.yorku.ca/SCS/vcd/)

**Discrete distributions**

<b>DISTPLOT</b>	Plots for discrete distributions
<b>GOODFIT</b>	Goodness-of-fit for discrete distributions
<b>ORDPLOT</b>	Ord plot for discrete distributions
<b>POISPLOT</b>	Poissonness plot
<b>ROOTGRAM</b>	Hanging rootograms

**Two-way and n-way tables**

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

**Model-based methods**

<b>ADDVAR</b>	Added variable plots for logistic regression
<b>CATPLOT</b>	Plot results from PROC CATMOD
<b>HALFNORM</b>	Half-normal plots for generalized linear models
<b>INFLGLIM</b>	Influence plots for generalized linear models
<b>INFLOGIS</b>	Influence plots for logistic regression
<b>LOGODDS</b>	Plot empirical logits and probabilities for binary data
<b>POWERLOG</b>	Power calculations for logistic regression
<b>POWERRxC</b>	Power calculations for two-way frequency table
<b>POWER2x2</b>	Power calculations for a $2 \times 2$ table
<b>ROBUST</b>	Robust fitting for linear models

**Utility macros**

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

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

**R software and the vcd package****Model-based methods**

<b>glm</b>	Fitting generalized linear models
<b>loglm</b>	MASS package: Fitting loglinear models
<b>R Commander</b>	Menu-driven package for statistical analysis and graphics
<b>car</b>	Package for graphics and extensions of generalized linear models
<b>effects</b>	Effects plots for generalized linear models

**R software and the vcd package**

- R software and the vcd package, available at [www.r-project.org](http://www.r-project.org)

**Discrete distributions**

<b>distplot</b>	Plots for discrete distributions
<b>goodfit</b>	Goodness-of-fit for discrete distributions
<b>ordplot</b>	Ord plot for discrete distributions
<b>poisplot</b>	Poissonness plot
<b>rootgram</b>	Hanging rootograms

**Two-way and n-way tables**

<b>agreementplot</b>	Observer agreement chart
<b>fourfold</b>	Fourfold displays for $2 \times 2 \times k$ tables
<b>sieveplot</b>	Sieve diagrams
<b>mosaic</b>	Mosaic displays
<b>pairs.table</b>	Matrix of pairwise association displays
<b>structable</b>	Manipulate high-dimensional contingency tables
<b>triplot</b>	Trilinear plots for $n \times 3$ tables

**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

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

$k$	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

## 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( $\lambda$ ) distribution.
- *Families in Saxony*: might expect a Bin( $n, p$ ) distribution with  $n = 12$ . Perhaps  $p = 0.5$  as well.

## Sidebar: Using SAS macros

- SAS macros are high-level, general programs consisting of a series of DATA steps and PROC steps.
- Keyword arguments substitute your data names, variable names, and options for the named macro parameters.
- Use as:
 

```
%macname(data=dataset, var=variables, ...);
```
- Most arguments have default values (e.g., data=\_last\_)
- All VCD macros have internal and online documentation, <http://www.math.yorku.ca/SCS/vcd/>
- Use as:
 

```
%macname(data=dataset, var=variables, ...);
```
- Macros can be installed in directories automatically searched by SAS. Put the following options statement in your AUTOEXEC.SAS file:
 

```
options sasautos=('c:\sasuser\macros' sasautos);
```

## 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 (**POISPLOT** macro)
- **Robust distribution plots** (**DISTPLOT** macro)

## Sidebar: Using SAS macros

E.g., the **GOODFIT** macro is defined with the following arguments:

```

... goodfit.sas ...
1 %macro goodfit(
2   data=_last_,      /* name of the input data set      */
3   var=,             /* analysis variable (basic count) */
4   freq=,            /* frequency variable              */
5   dist=,            /* name of distribution to be fit   */
6   parm=,            /* required distribution parameters? */
7   sumat=100000,     /* sum probs. and fitted values here */
8   format=,          /* format for ungrouped analysis variable */
9   out=fit,           /* output fit data set             */
10  outstat=stats); /* output statistics data set      */
```

Typical use:

```

1 %goodfit(data=madison,
2   var=count,
3   freq=blocks,
4   dist=poisson);
```

### Fitting discrete distributions

#### Distributions:

- Poisson,  $p(k) = e^{-\lambda} \lambda^k / k!$
- Binomial,  $p(k) = \binom{n}{k} p^k (1-p)^{n-k}$
- Negative binomial,  $p(k) = \binom{n+k-1}{k} p^n (1-p)^k$
- Geometric,  $p(k) = p(1-p)^k$
- Logarithmic series,  $p(k) = \theta^k / [-k \log(1-\theta)]$

#### Estimate parameter(s):

- Poisson,  $\hat{\lambda} = \sum k n_k / \sum n_k = \text{mean}$
- Binomial,  $\hat{p} = \sum k n_k / (n \sum n_k) = \text{mean} / n$

#### Goodness of fit:

$$\chi^2 = \sum_{k=1}^K \frac{(n_k - N \hat{p}_k)^2}{N \hat{p}_k} \sim \chi^2(K-1)$$

where  $\hat{p}_k$  is the estimated probability of each basic count, under the hypothesis that the data follows the chosen distribution.

### Fitting discrete distributions

The **GOODFIT** macro gives a table of observed and fitted frequencies, Pearson  $\chi^2$  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		

### 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);

```

### 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

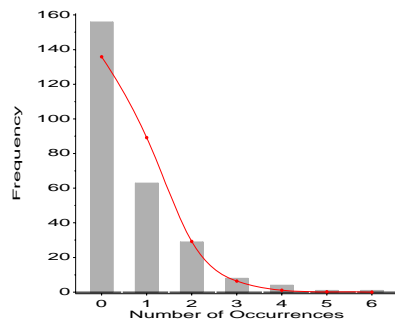
```

The poisson model does not fit! Why?

### 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);
```



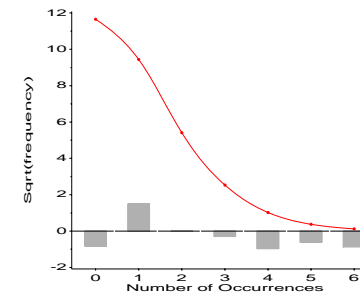
Problems:

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

### Highlight differences → Deviation rootograms

- Emphasize differences between observed and fitted frequencies
- Draw bars to show the gaps (btype=dev)

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

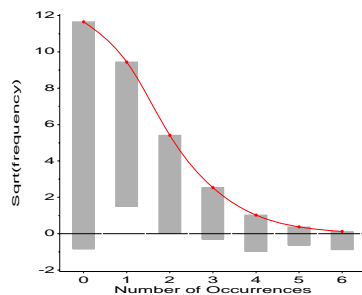


### Hang & root them → Hanging rootograms

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);
```



### 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  $kp_k/p_{k-1}$  against  $k$  is linear
  - signs of intercept and slope → determine the form, give rough estimates of parameters

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

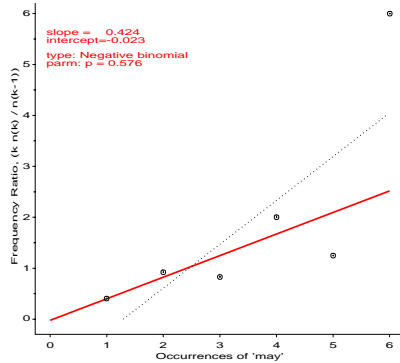
- Fit line by WLS, using  $\sqrt{n_k}$  as weights

### Ord plots

#### ■ ORD PLOT macro

```
%ordplot(data=madison, count=Count, freq=blocks);
```

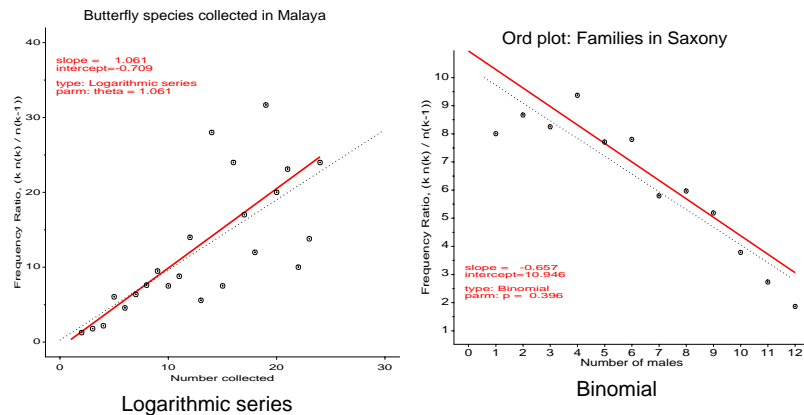
- Diagnoses distribution as NegBin
- Estimates  $\hat{p} = 0.576$



### 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 metamer**  $= \phi(n_k) = \log_e(k! n_k / N)$  vs.  $k$
  - Linear relation  $\Rightarrow$  Poisson, slope gives  $\hat{\lambda}$
  - CI for points, diagnostic (influence) plot
  - **POIS PLOT** macro

### Ord plots: Other distributions



### Poissonness plots: Details

- If the distribution of  $n_k$  is Poisson( $\lambda$ ) for some fixed  $\lambda$ , then each observed frequency,  $n_k \approx m_k = Np_k$ .
- Then, setting  $n_k = Np_k = e^{-\lambda} \lambda^k / k!$ , and taking logs of both sides gives

$$\log(n_k) = \log N - \lambda + k \log \lambda - \log k!$$

which can be rearranged to

$$\phi(n_k) \equiv \log\left(\frac{k! n_k}{N}\right) = -\lambda + (\log \lambda) k$$

- $\Rightarrow$  if the distribution is Poisson, plotting  $\phi(n_k)$  vs.  $k$  should give a line with
  - intercept =  $-\lambda$
  - slope =  $\log \lambda$
- Nonlinear relation  $\rightarrow$  distribution is *not* Poisson
- Hoaglin and Tukey (1985) give details on calculation of confidence intervals and influence measures.



### POISPLOT 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 %poispplot(data=madison,count=count, freq=blocks);

```

### Generalized robust distribution plots

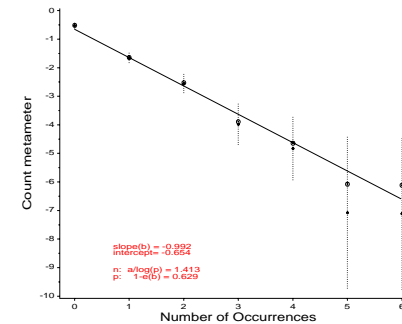
Other distributions: Analogous plots, for suitable count metameter,  $\phi(n_k)$  vs.  $k$ .

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

```

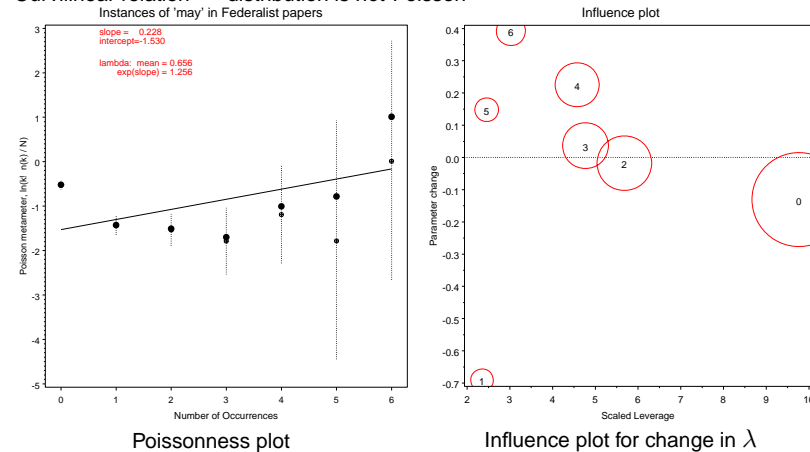
%distplot(data=madison, count=count, freq=blocks,
          dist=negbin);

```



### POISPLOT macro: output

Curvilinear relation  $\rightarrow$  distribution is *not* Poisson



### Testing Association in Two-Way Tables

#### ■ Typical analysis: Nominal factors

- Pearson  $\chi^2$  (or LR  $\chi^2$ )—when most expected frequencies  $\geq 5$ .

```

proc freq;
  weight count; /* if already a frequency table */
  table factor * response / chisq;

```

- Exact tests—small tables, small sample sizes (e.g., Fisher's)

```

proc freq;
  weight count; /* if already a frequency table */
  table factor * response / chisq;
  exact pchi;

```

### Example: Cholesterol and Heart disease

```

1 title 'Cholesterol diet and heart disease';
2 data fat;
3   input diet $ disease $ count;
4 datalines;
5 LoChol No 6
6 LoChol Yes 2
7 HiChol No 4
8 HiChol Yes 11
9 ;
10
11 proc freq data=fat;
12   weight count;
13   tables diet * disease / chisq nopercnt nocol;
14   exact pchi;

```

- Exact tests are valid and significant.

Exact test output:

```

-----
Pearson Chi-Square Test
-----
Chi-Square          4.9597
DF                  1
Asymptotic Pr > ChiSq 0.0259
Exact Pr >= ChiSq   0.0393

```

```

-----
Fisher's Exact Test
-----
Cell (1,1) Frequency (F) 4
Left-sided Pr <= F      0.0367
Right-sided Pr >= F     0.9967

Table Probability (P)    0.0334
Two-sided Pr <= P      0.0393

```

Standard output:

diet	disease		Total
	No	Yes	
HiChol	4	11	15
	26.67	73.33	
LoChol	6	2	8
	75.00	25.00	
Total	10	13	23

Statistics for Table of diet by disease

Statistic	DF	Value	Prob
Chi-Square	1	4.9597	0.0259
Likelihood Ratio Chi-Square	1	5.0975	0.0240
Continuity Adj. Chi-Square	1	3.1879	0.0742

WARNING: 50% of the cells have expected counts less than 5.  
(Asymptotic) Chi-Square may not be a valid test.

- The Pearson and LR  $\chi^2$  tests are not valid
- The conservative continuity-adjusted test fails significance

### Ordinal factors and Stratified analyses

#### More powerful CMH tests

- When either the row (factor) or column (response) levels are ordered, more specific (CMH = Cochran - Mantel - Haentzel) tests which take order into account have greater power to detect ordered relations.

```

proc freq;
  weight count;
  table factor * response / chisq cmh;

```

#### Control for other background variables

- Stratified analysis tests the association between a main factor and response *within* levels of the control variable(s)
- Can also test for homogeneous association across strata

```

proc freq;
  weight count;
  table strata * factor * response / chisq cmh;

```

### Example: Arthritis treatment

Data on treatment for rheumatoid arthritis (Koch and Edwards, 1988)

- **Ordinal response:** none, some, or marked improvement
- **Factor:** active treatment vs. placebo
- **Strata:** Sex

Treatment	Sex	Outcome			Total
		None	Some	Marked	
Active	Female	6	5	16	27
	Male	7	2	5	14
Placebo	Female	19	7	6	32
	Male	10	0	1	11
Total		42	14	28	84

### Overall analysis, ignoring sex: Results (chisq option)

STATISTICS FOR TABLE OF TREAT BY IMPROVE			
Statistic	DF	Value	Prob
Chi-Square	2	13.055	0.001
Likelihood Ratio Chi-Square	2	13.530	0.001
Mantel-Haenszel Chi-Square	1	12.859	0.000
Phi Coefficient		0.394	
Contingency Coefficient		0.367	
Cramer's V		0.394	

### Cochran-Mantel-Haenszel tests: (cmh option)

SUMMARY STATISTICS FOR TREAT BY IMPROVE				
Cochran-Mantel-Haenszel Statistics (Based on Table Scores)				
Statistic	Alternative Hypothesis	DF	Value	Prob
1	Nonzero Correlation	1	12.859	0.000
2	Row Mean Scores Differ	1	12.859	0.000
3	General Association	2	12.900	0.002

### Overall analysis, ignoring sex:

```

1 title 'Arthritis Treatment: PROC FREQ Analysis';
2 data arth;
3   input sex$ treat$ @;
4   do improve = 'None ', 'Some', 'Marked';
5     input count @;
6     output;
7   end;
8 datalines;
9 Female Active 6 5 16
10 Female Placebo 19 7 6
11 Male Active 7 2 5
12 Male Placebo 10 0 1
13 ;
14 *-- Ignoring sex;
15 proc freq order=data;
16   weight count;
17   tables treat * improve / cmh chisq nocol nopercnt;
18   run;

```

#### Notes:

- PROC FREQ orders character variables alphabetically (i.e., 'Marked', 'None', 'Some') by default. To treat the IMPROVE variable as ordinal, use order=data on the PROC FREQ statement.
- The chisq option gives the usual  $\chi^2$  tests (Pearson, Fisher's, etc.). The cmh option requests the Cochran-Mantel-Haenszel tests for ordinal variables.

### CMH tests for ordinal variables

- **Non-zero correlation:** Use when *both* row and column variables are ordinal.
  - $CMH \chi^2 = (N - 1)r^2$ , assigning scores (1, 2, 3, ...)
  - most powerful for *linear* association
- **Row Mean Scores Differ:** Use when only column variable is ordinal
  - Analogous to the Kruskal-Wallis non-parametric test (ANOVA on rank scores)
  - Ordinal variable should be listed *last* in the TABLES statement
- **General Association:** Use when *both* row and column variables are nominal.
  - Similar to overall Pearson  $\chi^2$  and Likelihood Ratio  $\chi^2$ .

### Sample CMH Profiles

Only general association:

	b1	b2	b3	b4	b5	Total	Mean
a1	0	15	25	15	0	55	3.0
a2	5	20	5	20	5	55	3.0
a3	20	5	5	5	20	55	3.0
Total	25	40	35	40	25	165	

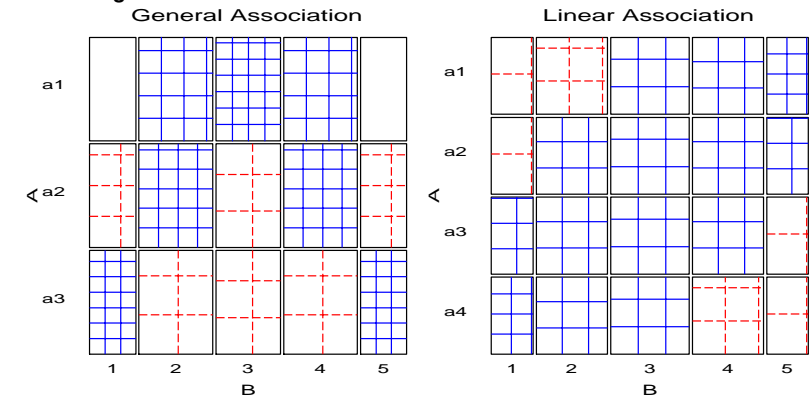
Output:

Cochran-Mantel-Haenszel Statistics (Based on Table Scores)

Statistic	Alternative Hypothesis	DF	Value	Prob
1	Nonzero Correlation	1	0.000	1.000
2	Row Mean Scores Differ	2	0.000	1.000
3	General Association	8	91.797	0.000

### Sample CMH Profiles

Visualizing Association:



### Sample CMH Profiles

Linear Association:

	b1	b2	b3	b4	b5	Total	Mean
a1	2	5	8	8	8	31	3.48
a2	2	8	8	8	5	31	3.19
a3	5	8	8	8	2	31	2.81
a4	8	8	8	5	2	31	2.52
Total	17	29	32	29	17	124	

Output

Cochran-Mantel-Haenszel Statistics (Based on Table Scores)

Statistic	Alternative Hypothesis	DF	Value	Prob
1	Nonzero Correlation	1	10.639	0.001
2	Row Mean Scores Differ	3	10.676	0.014
3	General Association	12	13.400	0.341

### Stratified analysis

#### Overall analysis

- ignores other variables (like sex), by collapsing over them
- risks losing important interactions (e.g., different associations for M and F)

#### Stratified analysis:

- controls for the effects of one or more background variables
- list stratification variable(s) *first* on the TABLES statement

```
tables age * sex * treat * improve;
```

### Stratified analysis

The statements below request a stratified analysis with CMH tests, controlling for sex.

... arthfreq.sas ...

```
20 *-- Stratified analysis, controlling for sex;
21 proc freq order=data;
22   weight count;
23   tables sex * treat * improve / cmh chisq nocol nopercnt;
24   run;
```

→ separate table (partial tests) for Females and Males

STATISTICS FOR TABLE 1 OF TREAT BY IMPROVE  
CONTROLLING FOR SEX=Female

Statistic	DF	Value	Prob
Chi-Square	2	11.296	0.004
Likelihood Ratio Chi-Square	2	11.731	0.003
Mantel-Haenszel Chi-Square	1	10.935	0.001
...			

- Strong association between TREAT and IMPROVE for females

### Stratified tests

- Individual (*partial*) tests are followed by a *conditional* test, controlling for strata (SEX)
- These tests **do not** require large sample size in the individual strata— just a large total sample size.
- They *assume*, but do not *test* that the association is the same for all strata.

SUMMARY STATISTICS FOR TREAT BY IMPROVE  
CONTROLLING FOR SEX

Cochran-Mantel-Haenszel Statistics (Based on Table Scores)

Statistic	Alternative Hypothesis	DF	Value	Prob
1	Nonzero Correlation	1	14.632	0.000
2	Row Mean Scores Differ	1	14.632	0.000
3	General Association	2	14.632	0.001

Males:

STATISTICS FOR TABLE 2 OF TREAT BY IMPROVE  
CONTROLLING FOR SEX=Male

Statistic	DF	Value	Prob
Chi-Square	2	4.907	0.086
Likelihood Ratio Chi-Square	2	5.855	0.054
Mantel-Haenszel Chi-Square	1	3.713	0.054
...			

WARNING: 67% of the cells have expected counts less than 5. Chi-Square may not be a valid test.

- Weak association between TREAT and IMPROVE for males
- Sample size  $N = 29$  for males is small

### Homogeneity of association

- is the association between the primary table variables is the same over all strata?
- $2 \times 2$  tables: → Equal odds ratios across all strata?
  - PROC FREQ: MEASURES option on TABLES statement → Breslow-Day test
- Larger tables: Use PROC CATMOD to test for *no three-way association* = same association for the primary factor & response variables.
- Arthritis data: homogeneity ↔ no 3-way sex \* treatment \* outcome association
  - ≡ loglinear model: [SexTreat] [SexOutcome] [TreatOutcome]
  - ≡ loglin sex|treat|improve@2 for PROC CATMOD
  - Zero frequencies: PROC CATMOD treats as “structural zeros” by default; recode if necessary.

... arthfreq.sas

```
26 title2 'Test homogeneity of treat*improve association';
27 data arth;
28   set arth;
29   if count=0 then count=1E-20;   *-- sampling zeros;
30 proc catmod order=data;
31   weight count;
32   model sex * treat * improve = _response_ / ml ;
33   loglin sex|treat|improve@2 / title='No 3-way association';
34 run;
35   loglin sex treat|improve / title='No Sex Associations';
```

### Homogeneity of association

- the likelihood ratio  $\chi^2$  (the badness-of-fit for the No 3-Way model) is the test for homogeneity
- clearly non-significant  $\rightarrow$  treatment-outcome association can be considered to be the same for men and women.

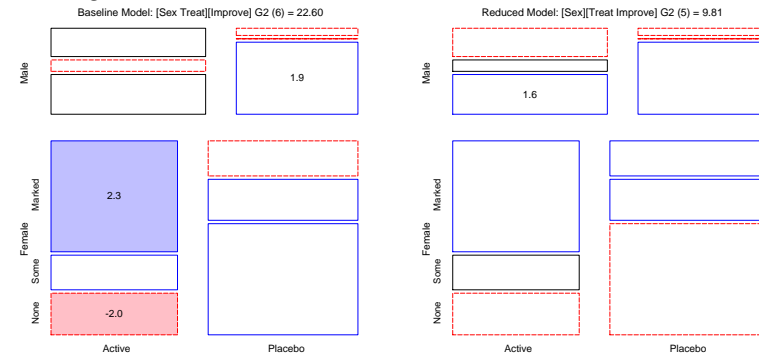
Test homogeneity of treat\*improve association  
No 3-way association  
MAXIMUM-LIKELIHOOD ANALYSIS-OF-VARIANCE TABLE

Source	DF	Chi-Square	Prob
SEX	1	14.13	0.0002
TREAT	1	1.32	0.2512
SEX*TREAT	1	2.93	0.0871
IMPROVE	2	13.61	0.0011
SEX*IMPROVE	2	6.51	0.0386
TREAT*IMPROVE	2	13.36	0.0013
LIKELIHOOD RATIO	2	1.70	0.4267

- But, associations of SEX\*TREAT and SEX\*IMPROVE are both small.
- Suggests stronger model of homogeneity, [Sex] [TreatOutcome], tested by `loglin sex treat|improve;` statement.

### Homogeneity of association

#### Visualizing Association:



### Homogeneity of association: Reduced model

```

30 proc catmod order=data;
31   weight count;
32   model sex * treat * improve = _response_ / ml ;
33   loglin sex|treat|improve@2 / title='No 3-way association';
34 run;
35   loglin sex treat|improve / title='No Sex Associations';

```

Output:

No Sex Associations  
MAXIMUM-LIKELIHOOD ANALYSIS-OF-VARIANCE TABLE

Source	DF	Chi-Square	Prob
SEX	1	12.95	0.0003
TREAT	1	0.15	0.6991
IMPROVE	2	10.99	0.0041
TREAT*IMPROVE	2	12.00	0.0025
LIKELIHOOD RATIO	5	9.81	0.0809

- Fits reasonably well
- How to interpret?