



Michael Friendly Psych 6136 ttp://friendly.github.io/613





# My goals

Start with descriptive, hypothesis testing methods, then progress to model-based methods

N

Sieve plots, mosaic plots, spineplots,

Visual tools for thinking & " understanding summ

Correspondence analysis: best 2D summary Effect plots, Data + Model plots

Build from simple, loglinear models to more complex ones

## 01: Overview

- Categorical data involves some new ideas
  - Discrete variables: unordered or ordered
  - Counts, frequencies as outcomes
- New / different data structures & functions
  - tables 1-way, 2-way, 3-way, ... table(), xtabs()
  - similar in matrices or arrays matrix(), array()
  - datasets:
    - frequency form
    - case form
- Graphical methods: often use area ~ Freq
  - Consider: graphical comparisons, effect order
- Models: Most are ≅ natural extensions of Im()

### Categorical data: Structures





The pattern is clearer when the eye colors are permuted: light hair goes with light eyes & vice-versa



1-way tables: graphs



Common discrete distributions

 $p^k(1-p)^{1-k}$ 

Success in 1

# successes

1st success

kth success

# of trials to 0.1.2

# of trials to 0, 1

in n trials

k={0, 1}

E(X)

1-p

. . . .

inference

p(1-p)

np(1-p)

1-p

k(1-p)

## 02: Discrete distributions

- Discrete distributions are the building blocks for categorical data analysis
  - Typically consist of basic counts of occurrences, with varying frequencies
  - Most common: binomial, Poisson, negative binomial
  - Others: geometric, log-series
- Fit with goodfit(); plot with rootogram()
  - Diagnostic plots: Ord\_plot(), distplot()
- Models with predictors
  - Binomial → logistic regression
  - Poisson → poisson regression; logliner models
  - These are special cases of generalized linear models



Graphing

Rootograms

s 6 7 8 9 10 11 mber of males	12	Poisson(A) Log series(p)	# of events in interval # of types observed	0, 1, 2,	$\frac{\frac{\lambda^{k}e^{-\lambda}}{k!}}{\frac{p^{i}}{n \log(1-p)}}$	٨	•
discrete distr	ibutions		Or	d plot	s: Examp	les	
Ord plots	Robust distribution plots	Ord plots fo > Ord_plot(Sa: > Ord_plot(Fe	or the Saxony xony, main = "Fa deralist, main = " Fanilies in taxes	and Federal milies in Saxony 'Instances of 'm	ist data ", gp=gpar(cex=1), pc ay' in Federalist pape instance	h=16) rs", gp=gpar(ce s of may is federalist	x=1), pch=16
		10 - integrate of the statement of the s	857 92946 93		6 steps = 0.424 intercept = -0.02 5 steps = 0.424 spec reincreati 4 estimate prote = 0 0 steps = 0 1 steps = 0 0 steps = 0 1 steps = 0 0 steps = 0 0 steps = 0 1 steps = 0 0 st	576 ••••••	•

## 03: Two-way tables

- Two-way tables summarize frequencies of two categorical factors
  - 2 × 2: a special case, with odds ratio as a measure
  - r × c: factors can be unordered or ordered
  - r × c × k: stratified tables, r × c with groups or circumstances
- Tests & measures of association
  - Pearson χ<sup>2</sup>, LR G<sup>2</sup>: general association
  - More powerful CMH tests for ordered factors
- Visualization
  - 2 × 2: fourfold plots
  - r × c: sieve diagrams, tile plots, ...
  - More graphical methods to come ...

### Measures of association





### CMH tests for ordinal factors

Pearson C

#### Three types of CMH tests:

Pearson contingency coef

- Non-zero correlation
- Use when both row and column variables are ordinal
   CMH χ<sup>2</sup> = (N 1)r<sup>2</sup>, assigning scores (1, 2, 3, ...)
- most powerful for *linear* association
- Row/Col Mean Scores Differ
- Use when only one variable is ordinal
- Analogous to the Kruskal-Wallis non-parametric test (ANOVA on rank scores)

#### General Association

Use when both row and column variables are nomina
 Similar to overall Pearson χ<sup>2</sup> and Likelihood Ratio G<sup>2</sup>

	8
•	Inter-observer agreement often used as to assess reliability of a subjective classification or assessment procedure • -> square table, Rater 1 x Rater 2
•	Agreement vs. Association: Ratings can be strongly associated without strong agreement

Observer agreement

 Marginal homogeneity: Different frequencies of category use by raters affects measures of agreement

#### Measures of Agreements

Intraclass correlation: ANOVA framework—multiple raters!
 Ochen's k: compares the observed agreement, P<sub>e</sub> = ∑p<sub>i</sub>, to agreement expected by chance if the two observe's ratings were independent, P<sub>e</sub> = ∑p<sub>i</sub>, p<sub>i</sub>, p<sub>i</sub>.

## 04: Loglinar models, mosaic displays

- Mosaic plots use sequential splits to show marginal and conditional frequencies in an *n*-way table
  - Shading: sign and magnitude of residuals  $\rightarrow$  contributions to  $\chi^2$
  - Shows the pattern of association not accounted for
  - Permuting rows/cols often helps
- Loglinear models
  - Express associations with ANOVA-like interaction terms: A\*B, A\*C
    - Joint independence: [AB][C] ≡ A \* B + C
    - Conditional independence: [AC][BC]  $\equiv$  A  $\perp$  B | C
  - Fitting models ≅ "cleaning the mosaic"
  - Response models: include all associations among predictors
- Sequential / partial plots & models
  - Sequential: Decompose all associations: V<sub>1</sub>; V<sub>2</sub>|V<sub>1</sub>; V<sub>3</sub>|{V<sub>1</sub>, V<sub>2</sub>}, ...
  - Partial: Decompose conditional associations: [V<sub>1</sub>, V<sub>2</sub>] | V<sub>3</sub>= {a, b, ...}

## Loglinear models: Perspectives

Loglinear models grew up and developed from three different ideas and ways of thinking about notions of independence in frequency data

- Loglinear approach: analog of ANOVA; associations are ~ interactions
- glm() approach: analog of general regression model, for log(Freq), with Poisson dist<sup>n</sup> of errors
- Logit models: Loglinear, simplified for a binary response

## Model-based methods: Fitting & graphing



## 05: Correspondence analysis

- CA is an exploratory method designed to account for association (Pearson  $\chi^2$ ) in a small number of dimensions
  - Row and column scores provide an optimal scaling of the category levels
  - Plots of these can suggest an explanation for association
- CA uses the singular value decomposition to approximate the matrix of residuals from independence
- Standard and principal coordinates have different geometric properties, but are essentially re-scalings of each other
- Multi-way tables can be handled by:
  - Stacking approach— collapse some dimensions interactively to a 2-way table
  - Each way of stacking  $\rightarrow$  a loglinear model
  - MCA analyzes the full n way table using an indicator matrix or the Burt matrix

Given a new 2-way table, my first thought is nearly always: plot(ca(mytable))

200	liced	model	c
1CU	uccu	noue	3

- For a three-way table there is a range of models between mutual independence, [A][B][C], and the saturated model, [ABC]
- Each model has an independence interpretation:
   [A][B] ≡ A ⊥ B ≡ A independent of B
- Special names for various submodels

Table: Log-linear	Models	for	Three-Way	Table

Model	Model symbol	Interpretation
Mutual independence	[A][B][C]	$A \perp B \perp C$
Joint independence	[AB][C]	$(A B) \perp C$
Conditional independence	[AC][BC]	$(A \perp B) \mid C$
All two-way associations	ABIACI(BC)	homogeneous assoc.
Saturated model	[ABC]	ABC interaction







- Rough interpretation: row/col points "near" each other are positively associated (independence residuals d<sub>2</sub> >> 0)
- Dim 1: 89.4% of χ<sup>2</sup> (dark → light)
   Dim 2: 9.5% of γ<sup>2</sup> (Red/Green vs. others)

### Multi-way tables: Stacking

A 3-way table of size  $I \times J \times K$  can be sliced and stacked as a two-way table in several ways



### Singular value decomposition

The singular value decomposition (SVD) is a basic technique for factoring a matrix and for matrix approximation For an  $n \times n$  matrix **X** of rank  $r \le \min(m, n)$  the SVD of **X** is:



## Multiple correspondence analysis

- Extends CA to n-way tables
- Useful when simpler stacking approach doesn't work well, e.g., 10 categorical attitude items
- Analyzes all pairwise bivariate associations. Analogous to:
   Correlation matrix (numbers)
  - Scatterplot matrix (graphs)
     All pairwise χ<sup>2</sup> tests (numbers)
     Mosaic matrix (graphs)
- Provides an optimal scaling of the category scores for each variable
- Can plot all factors in a single plot
- An extension, joint correspondence analysis, gives a better account of inertia for each dimension

## 06: Logistic regression

- loglm() provides only overall tests of model fit
- Model-based methods, glm(), provide hypothesis tests, CIs & tests for individual terms
- Logistic regression: A glm() for a binary response
  - Inear model for the log odds Pr(Y=1)
  - All similar to classical ANOVA, regression models
- Plotting
  - Conditional, full-model plots show data and fits
  - Effect plots show predicted effects averaged over others
- Model diagnostics
  - Influence plots are often informative

#### Modeling approaches: Overview Association models Response models Loglinear models Binary response (contingency table form) Categorical predictors: logit models [Admit][Gender Dept] logit(Admit) ~ 1 [Admit Dept][Gender Dept] logit(Admit) ~ Dep [AdmitDept][AdmitGender][Ge nderDept logit(Admit) ~ Dept + Gende Continuous/mixed predictors Poisson GLMs (Frequency data frame) Logistic regression models Freq ~ Admit + Gender \* Dept Pr(Admit) ~ Dent + Gender + Age + GRE Freq ~ Admit\*Dept + Gender\*Dept Freq ~ Admit\*(Dept + Gender) Polytomous response Gender\*Dent Ordinal: proportional odds mode Improve ~ Age + Sex + Treatment Ordinal variables General multinomial model Freq ~ right + left + Diag(right, left) WomenWork ~ Kids + HusbandIncome Freq ~ right + left + Symm(right, left) Full-model plot Plotting on the logit scale shows the additive effects of age, treatment and sex NB: easier to compare the treatment groups within the same panel 245.1 240 Treatmen Placebo

30 40 50 60 70 30 40 50 60 Age These plots show model uncertainty (confidence bands) Jittered points show the data



# 07: Logistic regression: Extensions

- Polytomous responses
  - *m* response categories  $\rightarrow$  (*m*-1) comparisons (logits)
  - Different models for ordered vs. unordered categories
- Proportional odds model
  - Simplest approach for ordered categories
  - Assumes same slopes for all logits
    - Fit with MASS::polr()
    - Test PO assumption with VGAM::vglm()
- Nested dichotomies
  - Applies to ordered or unordered categories
  - Fit m 1 separate independent models  $\rightarrow$  Additive G<sup>2</sup> values
- Multinomial logistic regression
  - Fit m 1 logits as a single model
  - Results usually comparable to nested dichotomies, but diff interpretation
  - R: nnet::multinom()



logit scale i.e.,  $\beta_1 = \beta_2$ .

Exploratory plots



When respons categories are

Ordered



Unordered

# **08: Extending loglinear models**

- Loglinear models, as originally formulated, were quite general, but treated all table variables as unordered factors
  - The GLM perspective is more general, allowing quantitative predictors and handling ordinal factors
  - The logit model give a simplified approach when one variable is a response
- Models for ordered factors give more powerful & focused tests
  - L × L, R, C and R+C models assign scores to the factors
  - RC(1) and RC(2) models estimate the scores from the data
- Models for square tables allow testing structured questions •
  - Quasi-independence: ignoring diagonals
  - symmetry & quasi-symmetry
  - theory-specific "topological" models
- These methods can be readily combined to analyze complex tables

### Logit models

40 60 80 10

For a binary response, each loglinear model is equivalent to a logit model (logistic regression, with categorical predictors) e.g., Admit ⊥ Gender | Dept (conditional independence = [AD][DG])

40 60 80 10

 $\log m_{ik} = \mu + \lambda_i^A + \lambda_i^D + \lambda_k^G + \lambda_{ii}^{AD} + \lambda_{ik}^{DG}$ So, for admitted (i = 1) and rejected (i = 2), we have  $\log m_{1jk} = \mu + \lambda_1^A + \lambda_j^D + \lambda_k^G + \lambda_{1j}^{AD} + \lambda_{ik}^{DC}$ 

(1)  $\log m_{2jk} = \mu + \lambda_2^A + \lambda_j^D + \lambda_k^G + \lambda_{2j}^{AD} + \lambda_{ik}^{DG}$ (2)

- Thus, subtracting (1)-(2), terms not involving Admit will cancel =  $\log m_{1/k} - \log m_{2/k} = \log(m_{1/k}/m_{2/k}) = \log \text{ odds of admission}$ 
  - $= (\lambda_1^A \lambda_2^A) + (\lambda_{1i}^{AD} \lambda_{2i}^{AD})$
  - $= \alpha + \beta_i^{\text{Dept}}$ (renaming terms)

where,  $\alpha$ : overall log odds of admission;  $\beta_i^{\text{Dept}}$ : effect on admissions of

## Square tables

Square tables arise when the row and column variables have the same categories, often ordered Special loglinear models allow us to tease apart different reasons for association



### Models for ordered categories

Consider an  $R \times C$  table having ordered categories

- In many cases, the RC association may be described more simply by assigning numeric scores to the row & column categories.
- For simplicity, we consider only integer scores, 1, 2, ... here

<ul> <li>These models are easily extended to stratified table</li> </ul>
--

R:C model	μ <sup>RC</sup>	df	Formula
Uniform association	$i \times j \times \gamma$	1	i:j
Row effects	$a_i \times j$	(I - 1)	R:j
Col effects	$i \times b_i$	(J - 1)	i:C
Row+Col eff	$ja_i + ib_i$	$\dot{I} + J - 3$	R:j + i:C
RC(1)	$\phi_I \psi_I \times \gamma$	I + J - 3	Mult(R, C)
Unstructured (R:C)	µRC .	(I-1)(J-1)	R:C

Model comparison plots When there are more than a few models, a model comparison plot can show the trade-off between goodness-of-fit and parsimony This sorts the models by both fit & complexity Plot BIC vs. df



Degrees

## 09: GLMs for Count Data

- GLMs provide a unified framework for linear models
  - Different families, all estimated in the same way
  - $\rightarrow$  link function and associated variance function
- For count data, starting from  $\log(\mu) = \mathbf{X} \beta$ ,  $\mu | \mathbf{X} \sim$ Poisson:
  - Overdispersion  $\rightarrow$  quasi-poisson, negative binomial
  - Standard tools for assessing model fit
- Excess zero counts introduce new ideas & methods ۲
  - ZIP model: structural model for the 0s
  - Hurdle model: random model for 0s, 2<sup>nd</sup> model for Y>0
- In all this, we rely on data & model plots for understanding

### Canonical links and variance functions

 For every distribution family, there is a default, canonical link function
 Each one also specifies the expected relation between the mean and variance

Table 11.2: Common di	istributions in the	exponential family	used with	generalized linear	models

ind their canonical link and variance functions						
Family	Notation	Canonical link	Range of y	Variance function, $\mathcal{V}(\mu \mid \eta)$		
Gaussian	$N(\mu, \sigma^2)$	identity: $\mu$	$(-\infty, +\infty)$	ø		
Poisson	$Pois(\mu)$	$\log_e(\mu)$	$0, 1, \ldots, \infty$	μ		
Negative-Binomial	$NBin(\mu, \theta)$	$\log_e(\mu)$	$0, 1,, \infty$	$\mu + \mu^2/\theta$		
Binomial	$\operatorname{Bin}(n,\mu)/n$	$logit(\mu)$	$\{0, 1,, n\}/n$	$\mu(1-\mu)/n$		
Gamma	$G(\mu, \nu)$	$\mu^{-1}$	$(0, +\infty)$	$\phi \mu^2$		
Inverse-Gaussian	$IG(\mu, \nu)$	$\mu^2$	$(0, +\infty)$	$\phi \mu^3$		

Choose a basic family:

- Get a default, canonical link, g(μ)
- Also get a variance function for free

#### Quasi-poisson models

- The quasi-<u>poisson</u> model allows the dispersion, φ, to be a free parameter, estimates with other coefficients
- The conditional variance is allowed to be a multiple of the mean

 $Var(y_i | \eta_i) = \phi \mu_i$ 

- This model is fit with glm() using family=quasipoisson
- The standard errors are multiplied by  $\phi^{\mbox{\tiny 22}}$
- Peace, order & good government is restored!

#### First, look at rootograms:

plot(goodfit(PhdPubs\$articles), xlab = "Number of Articles", main = "Polsson") plot(goodfit(PhdPubsSarticles, type = "nbinomial"), xlab = "Number of Articles", main = "Negative binomial")



One reason the Poisson doesn't fit: excess 0s (some never published?)

Q: What might some other reasons be? Think back to assumptions: independent obs; constant probs; unmodelled var:

### Models for excess zeros

- Two types of models, with different mechanisms for zero counts
- zero-inflated models: The responses with  $y_i = 0$  arise from a mixture of structural, always 0 values, with  $\Pr(y_i = 0) = \pi_i$  and the rest, which are random 0s, with  $\Pr(y_i = 0) = 1 \pi_i$
- hurdle models: One process determines whether  $y_i = 0$  with  $\Pr(y_i = 0) = \pi_i$ . A second process determines the distribution of values of positive counts,  $\Pr(y_i | y_i > 0)$

Zero-inflated



## 10: Models for log odds & LORs

- Logit models for a binary response generalize readily to a polytomous response
  - →Models for log odds, familiar interpretation
  - Handles 3+ way table, ordinal variables
  - Simple plots for interpretation
- Generalized odds ratios handle bivariate responses
  - Simple linear models for LOR
  - Easy to model log odds for each response and the LOR simultaneously
  - Easy to visualize results



# Your turn: Feedback?

## What did you like/dislike about 6136?

- Topics: what were the:
  - most interesting?
  - most boring?
  - Most challenging?
- What did you learn most from?
- What gave you the most difficulty?
- How does this relate to your own work?

# Tips for next time ...

- What should I try to differently the next time?
  - More of X?
  - Less of Y?
  - Aspects of how the course is structured?
  - Evaluation?

