

SFA 303 / 1002

Note Title

3/19/2012

TA Guest hours : today 2:00-4:00 RW 107/109
Friday 2:00-4:00 " "

Loglinear Models for Three-Dimensional Tables

Example: Alcohol, Cigarette & Marijuana Use
1992 survey of high-school seniors in Ohio
Have for each cell alcohol (A), cigarettes (C), marijuana (M)

DATA:

	<u>Alcohol Use</u>	<u>Cigarette Use</u>	<u>Marijuana Use</u>	
			<u>Yes</u>	<u>No</u>
<u>Yes</u>		<u>Yes</u>	911	538
		<u>No</u>	44	456
<u>No</u>		<u>Yes</u>	3	43
		<u>No</u>	2	279

Is there a response variable? Not a clear choice
Question of interest: Are variables associated?

General setup and notation: for 3-way contingency tables
Have data that are counts in $I \times J \times K$ table
Let π_{ijk} = probability an observation is classified
in cell (i, j, k)

Our example $I = J = K = 2$
 π_{ijk} denotes joint probability of
categorical random variables A, C, M
 y_{ijk} = observed count in cell (i, j, k)
 $\pi_{i..}$ / $y_{i..}$ = probability observation is in
level i of 1st variable /
count of observations in

level i of 1^{st} category
 $\pi_{ij} = \text{prob an observation is in level } i$
of 1^{st} variable and j of 2^{nd} var.
etc.

$$\pi_{i\dots} = 1$$

$$y_{i\dots} = n = \text{total of all covars}$$

y_{ijk} is a realization of a Poisson random variable
with mean μ_{ijk}
 $\hat{\mu}_{ijk} = n \pi_{ijk}$

Hierarchy of Models:

Model 1: Complete Independence

Short form: (A, C, M)

$$H_0: \pi_{ijk} = \pi_{i..} \pi_{.j.} \pi_{..k} \quad , \quad \begin{matrix} i=1, \dots, I \\ j=1, \dots, J \\ k=1, \dots, K \end{matrix}$$

$$H_a: \pi_{ijk} \neq \pi_{i..} \pi_{.j.} \pi_{..k}$$

If independent:

$$n \mu_{ijk} = n \pi_{ijk} = n \pi_{i..} \pi_{.j.} \pi_{..k}$$

$$\log(\mu_{ijk}) = \log n + \log \pi_{i..} + \log \pi_{.j.} + \log \pi_{..k}$$

Relative model

Alternative notation: $\log(\mu_{ijk}) = \alpha + \alpha_i^A + \alpha_j^C + \alpha_k^M$
(Agresti)

($1 = Y_{us}, 2 = N_{0}$)

Independence Model:

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbb{I}[A=1] + \beta_2 \mathbb{I}[C=1] + \beta_3 \mathbb{I}[M=1]$$

Find parameter estimates using maximum likelihood estimation,
subject to constraint $\sum_k \sum_i \sum_j \hat{\pi}_{ijk} = 1$

$$\Rightarrow \sum_k \sum_j \sum_i \hat{\mu}_{ijk} = n$$

$$= \sum_k \sum_j \sum_i y_{ijk}$$

Our show

$$\hat{\mu}_{ijk} = \left(n \hat{\pi}_{ijk} = n \hat{\pi}_{i..} \hat{\pi}_{.j.} \hat{\pi}_{..k} \right)$$

$$= n \frac{y_{i..}}{n} \frac{y_{.j.}}{n} \frac{y_{..k}}{n}$$

Estimates of μ 's make this true

Is the model assuming complete independence adequate?

Compare it to saturated model.

Saturated model includes all possible interactions:

Number of parameters in saturated model:

$$1 + 3 + 3 + 1 = 8$$

μ_0 $\mu_{A=1}, \mu_{A=2}, \dots, \mu_{M=1}$

1 2-way interactions 3-way interaction

In general:

$$1 + (I-1) + (J-1) + (K-1) + (I-1)(J-1) + (I-1)(K-1) + (J-1)(K-1) + (I-1)(J-1)(K-1) = IJK$$

\leftarrow Intercept + main effects + 2-way interactions + 3-way interactions

IJK table

= number of observed counts
Saturated model fits the data perfectly.

To compare fitted model (complete independence) to saturated model use $LRT = \text{Deviance goodness-of-fit}_{\text{test}}$

Example Deviance = 1286.62

Under H_0 , this is an observation from a chi square distribution with $g - k = 4$ df

Very strong evidence that saturated model fits better

Model 2!

Block independence

- add a 2-way interaction into the model
- 3 choices for which:
 - (A, C) - indicates model has main effect for M, A, C and interaction between A and C

or (A, M)

or (M, C)

The presence of a 2-way interaction between 2 variables indicates there is an association between these variables

log(LR, M) indicates alcohol & cigarette use are associated but they're probably independent of M

So - joint distribution of A, c, M factors into 2 blocks

For (A, c, M) :
 $H_0: \prod_{ijk} = \prod_{ij} \cdot \prod_{ik}$ $i=1, \dots, I$
 $H_a: \prod_{ijk} \neq \prod_{ij} \cdot \prod_{ik}$ $j=1, \dots, J$
 $k=1, \dots, K$

Model: (A, c, M)

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 I_{A=1} + \beta_2 I_{c=1} + \beta_3 I_{M=1} + \beta_4 I_{A=1} * I_{c=1}$$

Find est using ML estimation:

Can show: $\hat{\mu}_{ijk} (= n \hat{\pi}_{ijk} = n \hat{\pi}_{i \cdot \cdot k} \hat{\pi}_{\cdot j \cdot} \hat{\pi}_{\cdot \cdot k})$
 $= n \frac{y_{ij \cdot}}{n} \frac{y_{\cdot j \cdot}}{n}$

Is this an adequate model?

Use Deviance G-0-F test to compare to saturated model

Model (class) #3: Partial Independence

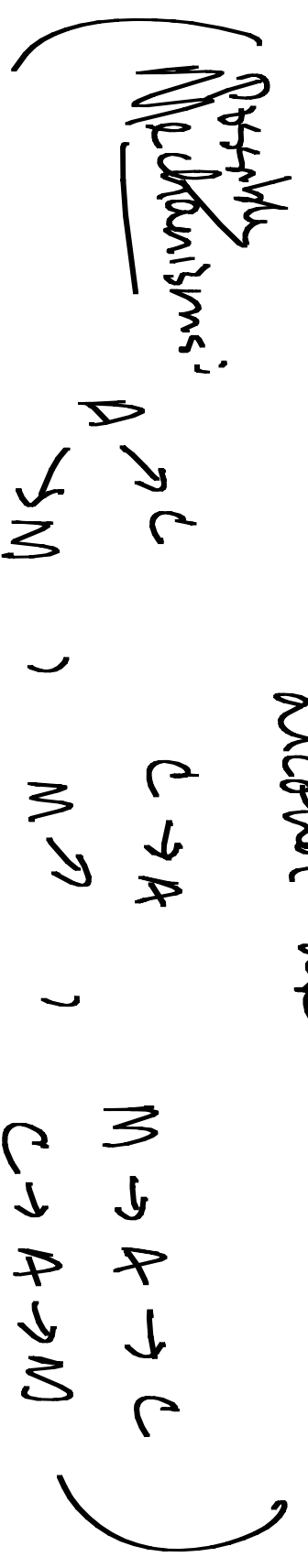
Model contains 2 of possible 2-way interactions,
so 2 pairs of variables are associated

Possibilities: (Ac, AM) *

or (Ac, CM)

or (AM, CM)

(Ac, AM) indicates alcohol & cigarette use are associated and alcohol and marijuana use are associated but cigarette and marijuana use are conditionally independent, conditional on alcohol use



Conditional independence:

$$P(CM | A) = P(C|A)P(M|A)$$

Translate to π 's

$$\frac{\pi_{ijk}}{\pi_{i..}} = \frac{\pi_{i.}}{\pi_{i..}} * \frac{\pi_{i.k}}{\pi_{i..}}$$

~~$P(y_1, y_2)$~~
 ~~$P(y_1)$~~
 ~~$P(y_2)$~~

or $\pi_{ijk} = \frac{\pi_{i.} \cdot \pi_{i.k}}{\pi_{i..}}$

MLEs given: $\hat{\mu}_{ijk} = n \hat{\pi}_{ijk} = n \frac{\hat{\pi}_{i.} \cdot \hat{\pi}_{i.k}}{\hat{\pi}_{i..}}$

(can ignore)

$$= n \frac{y_{i.}/n \cdot y_{i.k}/n}{y_{i..}/n}$$

$$= \frac{y_{ij} \cdot y_{i.k}}{y_{i..}}$$

Model 4: Uniform Association

(Ae, AM, cm)

- model contains all 2-way interactions
 - there is an association between each pair of variables, but no 3-way interaction
- \Rightarrow the way a pair of variables is associated is the same for all levels of the 3rd variable

For this model, There is no simple interpretation in terms of independence structure ; no nice form for μ_{ijk} 's
Need numerical procedures to solve for MLEs

Model 5: Saturated model
(ACM)
- perfectly fits table

Example

Model	G^2 (Deviance)	df	p
A, C, M	1286	4	< .0001
A, CM	534	3	< .0001
AM, C	940	3	< .0001
AC, M	844	3	< .0001
AC, AM	497	2	< .0001
AC, CM	92	2	< .0001
AM, CM	184	2	< .0001
AC, AM, CM	.4	1	<u>.54</u>

SAS doesn't give

About simplest model that fits data adequately:
uniform association

Key to likelihood model analysis: examine association
(NOT how some variables predict / explain a response)

Example: Uniform Association

Fitted equation:

$$\log(\hat{\mu}_{ijk}) = 5.63 + .049 I_A - 1.89 I_C - 5.31 I_M \\ + 2.05 I_A * I_C + 2.99 I_A * I_M + 2.85 I_C * I_M$$

Example calculation of a predicted value:

$$\begin{aligned}\hat{\mu}_{111} & (A=Yes, C=Yes, M=Yes, I_A=1, I_C=1, I_M=1) \\ & = \exp\{5.63 + .049 - 1.89 + \dots + 2.85\} = 910.38\end{aligned}$$

$$\begin{aligned}\hat{\mu}_{222} & (A=No, C=No, M=No, I_A=I_C=I_M=0) \\ & = \exp\{5.63\}\end{aligned}$$

Interpreting $\hat{\beta}$'s — use them to calculate odds

$$\hat{\mu}_{ijk} = N \hat{\pi}_{ijk}$$

gives odds of marijuana use for i, j
 alcohol & cigarette use = i, j

$$\frac{\hat{\pi}_{ij1}}{\hat{\pi}_{ij2}}$$

$$= \frac{\hat{N}_{ij1}}{\hat{N}_{ij2}}$$

Example

For students use alcohol & cigarettes, ^{estimated} odds of
 marijuana use are $\frac{910.38}{538.62}$

← from SATS predicted count for
 $A=1, C=1, M=2$
 $= 1.7$

For students who used neither alcohol or cigarettes,

estimated odds of marijuana use

$$\frac{1.38}{279.62}$$