

# STA 303/1002-Methods of Data Analysis II

Sections L0101& L0201, Winter 2018

**Shivon Sue-Chee**



March 15, 2018

Three-way Contingency Tables

## Class 18- Summary of Case Study VI



### Framingham Heart Study

A Project of the National Heart, Lung, and Blood Institute and Boston University

	$I$ Diff. in prop	LRT	Log-linear
Assume	Row totals fixed	Overall total fixed	Totals are random
Dist. of Y	Binomial	Multinomial	Poisson
$H_0$	$\pi_1 = \pi_2$	$\pi_{ij} = \pi_{i.} \pi_{.j}$	Additive model
Test Stat.	Z	$\chi^2_{(I-1)(J-1)}$	$\chi^2_{(I-1)(J-1)}$

$2 \times 2$  CT  
 $\downarrow$   
 $(2 \text{ factors}) \quad I \times J \text{ CT}$   
 $I, J > 2$   
 $\downarrow$   
 $(3 \text{ factors}) \quad I \times J \times K \text{ CT}$   
 Eg.  $2 \times 2 \times 2$   
 Non-parametric method: Fisher's Exact Test.

Three-way Contingency Tables  
 $\text{table } (f_1, f_2, f_3) \leftarrow \dim(f_3) \text{ 2-way tables}$

## A Three-way Contingency Table

Case Study VII Data:

- ▶ 1992 survey of high-school seniors in Ohio
- ▶ Table of counts of seniors who used alcohol, cigarettes and marijuana.

(3 factors).

Alcohol use	Cigarette use	Marijuana use	
		Yes	No
Yes	Yes	911	538
	No	44	456
No	Yes	3	43
	No	2	279

Q: Are alcohol (A), cigarettes (C) and marijuana (M) use associated?

## Forms of independence in $I \times J \times K$ Tables

Models

	Independence	$\pi_{ijk}$	Short form	
	Mutually indep.	(1) $\pi_{ijk} = \pi_{i..}\pi_{.j.}\pi_{..k}$	(X,Y,Z)	Complete
1 2-way	Jointly indep.	(3) $\pi_{ijk} = \pi_{ij.}\pi_{..k}$	( <span style="border: 1px solid red;">XY</span> ,Z)	Block
2 2-way	<u>Conditionally indep.</u>	(3) $\pi_{ijk} = \pi_{i.k}\pi_{.jk}/\pi_{..k}$	(XZ,YZ)	Partial
All 2-way	Uniform assoc.	(1) $\pi_{ijk} = \pi_{ij.}\pi_{i.k}\pi_{.jk}$	(XZ,YZ, XY)	Homo
All 2-way & 3-way	Saturated	(1) $\pi_{ijk}$	XYZ	

Eg

$\left\{ \begin{array}{l} (AC, M) \\ (AM, C) \\ (CM, A) \end{array} \right\}$   
 $\left\{ \begin{array}{l} (XY, XZ) \\ (YZ, XY) \end{array} \right\}$

Three-way Contingency Tables

# Three-way Tables



## ► Learning Objectives

- Write out the models used and the assumptions for inference
- Carry out the inference procedures completely
- Interpret the respective R outputs

$$\ln(\mu_{ijk}) = \beta_0 + \beta_1 I_A + \beta_2 I_C + \beta_3 I_M$$

$$- - - \frac{1}{A} \frac{1}{C}$$

Wald

G-O-F

Global LRT

$$- - - \frac{1}{A} \frac{1}{C} \frac{1}{M}$$

Three-way Contingency Tables

## Model 1: Complete Independence

- ▶  $P(ACM) = P(A)P(C)P(M)$ ; Alcohol, cigarette and marijuana use are **mutually independent**

- ▶ Hypotheses:

$$H_0 : \pi_{ijk} = \pi_{i..}\pi_{.j.}\pi_{..k} \text{ for all } i, j, k$$

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

- ▶ Short form: (A,C,M) -all 3 main effects only
- ▶  $I = J = K = 2$

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M$$

Additive

where  $\mathbf{I} : \{1 = \text{Yes}, 0 = \text{No}\}$

## Model 1: Complete Independence

- In general, we have the constraint  $n = \sum_i \sum_j \sum_k y_{ijk}$  or  $\sum_i \sum_j \sum_k \hat{\pi}_{ijk} = 1$ . Then by ML estimation,

$$\begin{aligned} \sum_i \sum_j \sum_k \hat{\mu}_{ijk} &= n = \sum_i \sum_j \sum_k y_{ijk} \\ \implies \hat{\pi}_{ijk} &= \frac{y_{ijk}}{n} \text{ or } \hat{\mu}_{ijk} = y_{ijk} \end{aligned}$$

- For complete independence model, using an additional  $(I - 1) + (J - 1) + (K - 1)$  constraints

$$\begin{aligned} \hat{\mu}_{ijk} &= n \hat{\pi}_{ijk} = n \hat{\pi}_{i..} \hat{\pi}_{.j.} \hat{\pi}_{..k} \\ &= n \frac{y_{i..}}{n} \frac{y_{.j.}}{n} \frac{y_{..k}}{n} \end{aligned}$$

LRT  
(Deviance G-0-F)

$H_0$ : Fitted  
 $H_a$ : Saturated

## Model Class 2: Block Independence

- ▶  $P(AC|M) = P(AC)$ ; Joint probability of alcohol and cigarette use is independent of marijuana use; Alcohol and cigarette use are associated
- ▶ Hypotheses:

$$H_0 : \pi_{ijk} = \pi_{ij.}\pi_{..k}$$

$$H_a : \pi_{ijk} \neq \pi_{ij.}\pi_{..k}$$

- ▶ Short form: (AC,M) - all 3 main effects and 1 interaction

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M + \beta_4 \mathbf{I}_{AC}$$

where  $\mathbf{I}_{AC} = \mathbf{I}_A * \mathbf{I}_C$

- ▶ Others in this class: (AM, C), (CM, A)

1 2-way interaction term.



## Model 2: Block Independence

- By ML estimation, for block independence model

$$\begin{aligned}\hat{\mu}_{ijk} &= n\hat{\pi}_{ijk} = n\hat{\pi}_{ij\cdot}\hat{\pi}_{\cdot\cdot k} \\ &= n\frac{y_{ij\cdot}}{n}\frac{y_{\cdot\cdot k}}{n}\end{aligned}$$

## Model Class 3: Partial Independence

- ▶  $P(AC|M) = P(A|M)P(C|M)$ ; Alcohol and cigarette use are conditionally independent given marijuana use; Alcohol and marijuana use are associated, and cigarette and marijuana use are associated
- ▶ Hypotheses:

$$H_0 : \pi_{ijk} = \pi_{i \cdot k} \pi_{\cdot j k} / \underline{\underline{\pi_{\cdot \cdot k}}}$$

$$H_a : \pi_{ijk} \neq \pi_{i \cdot k} \pi_{\cdot j k} / \underline{\underline{\pi_{\cdot \cdot k}}}$$

- ▶ Short form: (AM,CM) - all 3 main effects and 2 interactions

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M + \beta_4 \mathbf{I}_{AM} + \beta_5 \mathbf{I}_{CM}$$

- ▶ Others in this class: (AC, CM), (AC, AM)

2 2-way  
interaction  
terms.

## Model 3: Partial Independence

- ▶ We have  $P(AC|M) = P(A|M)P(C|M)$ .

$$\begin{aligned}\implies \frac{\pi_{ijk}}{\pi_{..k}} &= \frac{\pi_{.jk}}{\pi_{..k}} \frac{\pi_{i.k}}{\pi_{..k}} \\ \text{or } \pi_{ijk} &= \frac{\pi_{.jk}\pi_{i.k}}{\pi_{..k}}\end{aligned}$$

- ▶ Then by ML estimation

$$\begin{aligned}\hat{\mu}_{ijk} &= n\hat{\pi}_{ijk} = n \frac{\hat{\pi}_{.jk}\hat{\pi}_{i.k}}{\pi_{..k}} \\ &= \cancel{n} \frac{(y_{.jk}/n)(y_{i.k}/n)}{(y_{..k}/n)} \\ &= \frac{y_{.jk}y_{i.k}}{y_{..k}}\end{aligned}$$

## Model 4: Uniform association

- ▶ There is an association among all pairs
- ▶ For all levels of the 3rd variable, the association between the pair is the same
- ▶ Short form: (AM,AC,CM) - all 3 main effects and 3 two-way interactions but no three-way interaction

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M + \beta_4 \mathbf{I}_{AM} + \beta_5 \mathbf{I}_{AC} + \beta_6 \mathbf{I}_{CM}$$

All 2-way interaction terms.

- ▶ Solutions for  $\pi_{ijk}$  ( $\mu_{ijk}$ ) are found numerically with no simple expression in terms of  $y_{ijk}$ 's
- ▶ No simple interpretation to independence structure

## Saturated Model

- ▶ Total number of parameters:

$$1 + \underbrace{3}_{1\text{-way}} + \underbrace{3}_{2\text{-way}} + \underbrace{1}_{3\text{-way}} = 8$$

- ▶ Total number of observed counts:

$$\begin{aligned} &1 + (I - 1) + (J - 1) + (K - 1) \\ &+ (I - 1)(J - 1) + (I - 1)(K - 1) + (J - 1)(K - 1) \\ &+ (I - 1)(J - 1)(K - 1) = IJK = 2 * 2 * 2 = 8 \end{aligned}$$

$$\begin{aligned} \log(\mu_{ijk}) = &\beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M \\ &+ \beta_4 \mathbf{I}_{AM} + \beta_5 \mathbf{I}_{AC} + \beta_6 \mathbf{I}_{CM} + \beta_7 \mathbf{I}_{ACM} \end{aligned}$$

- ▶ Saturated model always fits the data perfectly

## On the Saturated Model

$$\log(\mu_{ijk}) = \beta_0 + \beta_1 \mathbf{I}_A + \beta_2 \mathbf{I}_C + \beta_3 \mathbf{I}_M \\ + \beta_4 \mathbf{I}_{AM} + \beta_5 \mathbf{I}_{AC} + \beta_6 \mathbf{I}_{CM} + \beta_7 \mathbf{I}_{ACM}$$

- ▶ Total # of parameters = Total # of observed counts
- ▶ Has a separate parameter for each observation
- ▶ Always gives a perfect fit
- ▶ Explains all the variation by its systematic component
- ▶ Sounds good but not a helpful model
- ▶ Does not smooth the data or is not parsimonious
- ▶ Serves as a baseline for checking model fit

Three-way Contingency Tables

Deviance G-O-F

$H_0$ : Fitted  
(Reduced)

$H_a$ : Saturated  
(Full).

$G-S = \text{Deviance}_R$   
—  
~~Deviance<sub>S</sub>~~ ↓ 0

## Results from R output

Use R and codes provided to show these

→ Add AIC statistics to the table

Model	df	$G^2$ =Deviance	p-value
(A,C,M)	4	1286.02	< 0.0001
(AC,M)	3	843.83	< 0.0001
(AM, C)	3	939.56	<0.0001
(A,CM)	3	534.21	<0.0001
(AC,AM)	2	497.37	<0.0001
(AC,CM)	2	92.02	<0.0001
(AM,CM)	2	187.75	<0.0001
(AC,AM,CM)	1	0.37	0.5408
(ACM)	0	0.00	-

fitted model is not adequate

Adequate.

The simplest model that fits the data adequately is the "Uniform Association" model (AC,AM,CM).



## Class 18 Summary

- ▶ Three-way contingency tables:

- ▶ Log-linear model approach
- ▶ Types of independence or association/ interactions
  - (i) Complete
  - (ii) Block
  - (iii) Partial
  - (iv) Uniform association
  - (v) 3-way interaction
- ▶ Deviance goodness-of-fit test

- ▶ Next Class: Model Diagnostics

*Interpretation, Estimate*

- ▶ Things to do:

- ▶ Assignment #3
- ▶ Participation 6
- ▶ Participation 7
- ▶ Practice Problems on Poisson Regression (Log-linear models)

Three-way Contingency Tables

*Final (Apr. 25)*