

STA 303 / 1002

Note Title

3/12/2012

Assignment 2 is ready!

The TA will be in RW 107/109 to help with the assignment on

Mon. March 19	2:00 - 4:00 pm
Fri. March 23	2:00 - 4:00 p.m.

Note: The guest server has moved. If you login remotely (ssh) you may get an error about speaking. Follow the instructions in the announcement on Blackboard.

Framingham Example

Status	CVD		Total
	present	absent	
High: ≥ 260	41	245	286 = N_H
Low: < 260	51	992	1043 = N_L
Total	92	1237	1329

$H_0: \pi_H = \pi_L$ where π_H is prob. of developing CVD in high cholesterol group

$H_a: \pi_H \neq \pi_L$

Analysis based on binomial distribution

~~Assumptions:~~ 286, 1043 are fixed,
i people independent

$$\hat{\pi}_H = 41/286, \quad \text{Var}(\hat{\pi}_H) = \frac{\pi_H(1-\pi_H)}{n_H}$$

$$\hat{\pi}_L = 51/1043, \quad \text{Var}(\hat{\pi}_L) = \frac{\pi_L(1-\pi_L)}{n_L}$$

Test statistic:

$$\frac{\hat{\pi}_H - \hat{\pi}_L}{\text{s.d. of } (\hat{\pi}_H - \hat{\pi}_L)}$$

$$\text{Var}(\hat{\pi}_H - \hat{\pi}_L) = \text{Var}(\hat{\pi}_H) + \text{Var}(\hat{\pi}_L)$$

Under H_0 , $\pi_H = \pi_L$ so estimate them by

$$\hat{\pi}_{\text{combined}} = 92/1329$$

$$\text{Under } H_0, \text{ S.E. of } (\hat{\pi}_A - \hat{\pi}_L) =$$

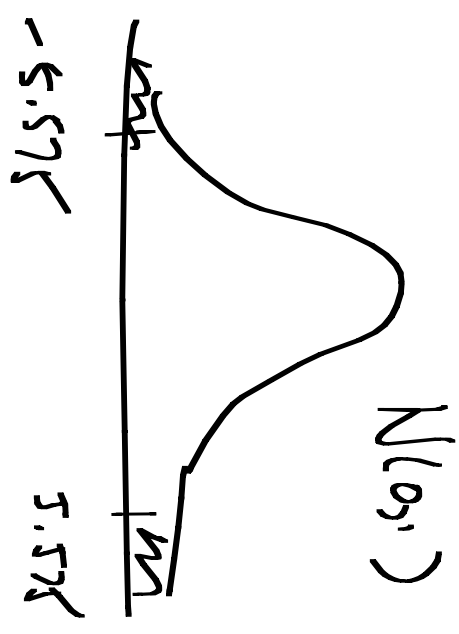
$$\sqrt{\frac{\frac{92/1329 (1 - 92/1329)}{286} + \frac{92/1329 (1 - 92/1329)}{1043}}$$

For large samples, proportions are approximately normally distributed (ANT)

So test statistic is approx N, dist'd under H_0

For these data, the test statistic is 5.575
The p-value is very small

Very strong evidence that
probability of developing CVD
is different for high vs low
cholesterol groups



Underlying probability distribution: Binomial

"Product binomial" or "Binomial sampling"
approach

ANALYSIS 2 Assume sample $n = 1329$ people
(concordance fixed),

classified them 2 ways

cholesterol status

(H or L)

A/D Status

(present or absent)

Let C denote cholesterol status
and D denote disease status

Both categorical
r.v.'s
(with 2 categories)

In general, row factor with I levels
column factor with J levels

π_{ij} = probability an observation falls into row i
column j
 $= P_r (C=i, D=j)$

gives joint distribution of C & D

Marginal distributions

$P(C=i) = \pi_{i.} =$ probability an observation falls into row i

$$P(D=j) = \pi_{.j} = \begin{matrix} n & n & n \\ \vdots & \vdots & \vdots \\ n & n & n \end{matrix} \text{ column } j$$

If there is no relationship between C and D ,
 C, D are independent

That is, $\pi_{ij} = \pi_{i.} \cdot \pi_{.j}$

$H_0: \pi_{ij} = \pi_{i.} \pi_{.j}, i=1, \dots, I, j=1, \dots, J$

vs $H_{a1}: \pi_{ij} \neq \pi_{i.} \pi_{.j}$

Let y_{ij} = count in row i , column j
 $y_{i.}$ = total count in row i , $y_{.j}$ = total count in column j
 y_{ij} = " " " " column j = $\sum_{i=1}^I y_{ij}$

$$n = \sum_i \sum_j y_{ij} = \text{grand total}$$

Under H_0 , estimate of expected count in cell (i, j)
 $\mu_{ij} = n \hat{\pi}_{i.} \hat{\pi}_{.j}$

$$\begin{aligned}
 &= n \frac{y_{i \cdot}}{n} \frac{y_{\cdot j}}{n} = \frac{y_{i \cdot} y_{\cdot j}}{n} \\
 \text{Test statistic} & \quad \chi^2 = \sum_{j=1}^J \sum_{i=1}^I \frac{(y_{ij} - \hat{\mu}_{ij})^2}{\hat{\mu}_{ij}}
 \end{aligned}$$

For large samples, under H_0 χ^2 has approximately a Chi-square distribution with $df = (I-1)(J-1)$

Why $(I-1)(J-1)$ df ?

IJ observations

Restrictions based on number of estimates we

How to calculate test statistic

Distribution on df

To test $\hat{\mu}_{ij}$, need $y_{i \cdot}$, $I-1$ estimates
 $y_{\cdot j}$, $J-1$ estimates
(I rows must sum to n)

n fixed.

$$df = IJ - 1 - (I-1) - (J-1) \\ = (I-1)(J-1)$$

SAS: proc freq; ("freq" short for frequency)
— For this analysis: Chi-square statistic

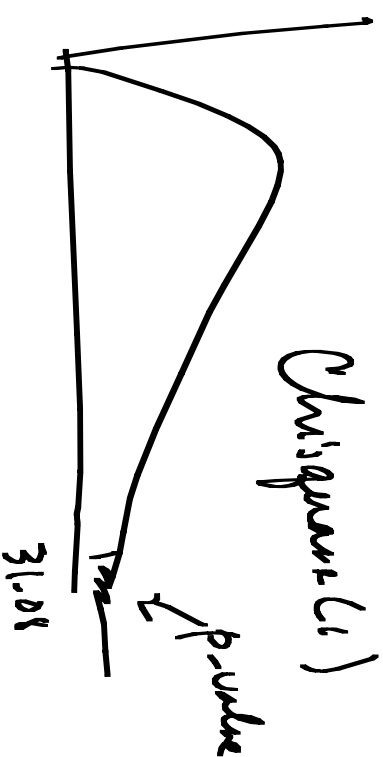
For example $df = 1$ ($J = I = 2$)

$$\chi^2 = 31.08$$

$$p\text{-value} < .0001$$

Strong evidence that C:ID are not independent

So C:ID status depends on cholesterol status



Exercise!

Show χ^2 equal to square of test statistic in Analysis! when $I = J = 2$

Analysis 2b

$n = 1329$ observations (fixed)
classified 2 ways

More formal approach based on MLE & LRTs:
Let y_i be a r.v., representing the number
of observations in row i , column j of table

observe y_i

For $I = J = 2$

Multinomial distribution:

$$P(Y = y) =$$

$$Y = (Y_{11}, Y_{12}, Y_{21}, Y_{22})$$
$$\frac{n!}{y_{11}! y_{12}! y_{21}! y_{22}!} \pi_{11}^{y_{11}} \pi_{12}^{y_{12}} \pi_{21}^{y_{21}} \pi_{22}^{y_{22}}$$

Understandability
per person
Multinomial

log-likelihood:

$$\log L = \sum_{j=1}^J \sum_{i=1}^I y_{ij} \log(\pi_{ij}) + \text{constant term}$$

$\underbrace{\hspace{10em}}$
 SAs would call log-likelihood
 $\underbrace{\hspace{10em}}$
 SAs would call full log-likelihood

$\underbrace{\hspace{10em}}$
 Number of ways of observations so
 distributed in row, col,
 that y_{11} are in row 1, col 1,
 etc

with $y_{11} + y_{12} + y_{21} + y_{22} = n$

and $\pi_{11} + \pi_{12} + \pi_{21} + \pi_{22} = 1$

Maximum likelihood estimates:

$$\hat{\pi}_{ij} = \frac{y_{ij}}{n}$$

(found by maximizing log L w.r.t. $\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22}$ subject to constraint $\sum \pi_{ij} = 1$)

Under assumption of independence:

$$\pi_{ij} = \pi_{i.} \pi_{.j}$$

Substitute $\pi_{ij} = \pi_{i.} \pi_{.j}$

into log L and maximize w.r.t $\pi_{1.}, \pi_{2.}, \pi_{.1}, \pi_{.2}$ subject to constraints $\pi_{1.} + \pi_{2.} = 1$ and $\pi_{.1} + \pi_{.2} = 1$

Solution

MLEs under assumption of independence:

$$\hat{\pi}_{1.} = \frac{y_{1.}}{n}, \quad \hat{\pi}_{2.} = \frac{y_{2.}}{n}$$

$$\hat{\pi}_{.1} = \frac{y_{.1}}{n}, \quad \hat{\pi}_{.2} = \frac{y_{.2}}{n}$$

Then $\hat{\pi}_{12} = \hat{\pi}_{1.} \hat{\pi}_{.2}$, etc., etc.

Leads to same expected counts as χ^2

Likelihood ratio test to compare multinomial model

under assumption of independence to model without this

assumption:

↖ "reduced model"

↖ "full model"

$$\text{Test } G^2 = -2 \log\left(\frac{L_E}{L_F}\right)$$

$$= 2 \log(L_E) - 2 \log(L_F) \leftarrow \hat{\pi}_{ij} \text{ for full model}$$

$$= 2 \left\{ \sum_i \sum_j y_{ij} \log\left(\frac{y_{ij}}{n}\right) - \sum_i \sum_j y_{ij} \log\left(\frac{y_{i\cdot} y_{\cdot j}}{n}\right) \right\}$$

$$= 2 \left\{ \sum_i \sum_j y_{ij} \log\left(\frac{y_{ij}}{n}\right) - \sum_i \sum_j y_{ij} \log\left(\frac{y_{i\cdot} y_{\cdot j}}{n}\right) \right\}$$

$\hat{\pi}_{ij}$ for reduced model

$$= 2 \sum_i \sum_j y_{ij} \log\left(\frac{y_{ij}}{\hat{\pi}_{ij}}\right)$$

where $\hat{\mu}_{ij}$ is expected count in cell i, j
 under assumption of independence

Under H_0 : $(H_0: \pi_{ij} = \pi_{i.} \pi_{.j} \text{ vs } H_a: \pi_{ij} \neq \pi_{i.} \pi_{.j}, i=1, \dots, I, j=1, \dots, J)$

G^2 has a chi-square distribution with $df = (I-1)(J-1)$

Why $(I-1)(J-1)$ df?

Unrestricted model: number of parameters $(\pi_{ij}'s)$ is $IJ - 1$ (use 1 df because $\sum \sum \pi_{ij} = 1$)

Reduced model (independence): number of parameters

(parameters are π_i, π_j 's)
(lose 2 df because
 $\sum_i \pi_i = 1, \sum_j \pi_j = 1$)
so # of parameters
is $I+J-2$

df for test statistic in LRT is

$$IJ-1 - (I+J-2) = (I-1)(J-1)$$

SS for freq, $G^2 =$ 'likelihood ratio chi-square'

For our example, $G^2 = 26.4298$
df = 1

P-value < 0.05

Strong evidence that row and column variable (cholesterol & disease status) are not independent

[Not responsible for numbers from pre-freq output]
[When likelihood ratio Chi-square]

What is the relationship between CVD & cholesterol?
For people with high cholesterol, 14.3% developed CVD
For people with low cholesterol, 4.99% developed CVD

Analysis 3 No fixed counts

Treat the IT counts as realizations of independent Poisson random variables.

The joint of counts is $P(\underline{y} = \underline{y}) = \prod_j \prod_i \frac{\mu_{ij}^{y_{ij}} e^{-\mu_{ij}}}{y_{ij}!}$

log-likelihood function:

$$\log L = \sum_j \sum_i (y_{ij} \log(\mu_{ij}) - \mu_{ij} - \log(y_{ij}!))$$

"full log-likelihood"

Use Poisson Regression (predefined in STAs)
" log-linear " models