

5) We want to test

$H_0$ : there's no relationship between response and ethnic origin

We use chi-square test of independency  $H_0: \pi_{ij} = \pi_i \cdot \pi_j$

$$X^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2_{(I-1)(J-1)} \quad \text{reject } H_0 \text{ if } X^2 \text{ too big}$$

where  $O_{ij}$ : observed frequencies

$I$ : cells = number of rows = 9

$E_{ij}$ : expected frequencies

$J$ : mult. dist = number of columns = 2

Origin	Yes	No	Total
Italian	78	47	125
Nuth. Eu	56	29	85
Other Eu	43	29	72
English	53	32	85
Irish	43	30	73
Fr. Canadian	36	22	58
French	42	23	65
Portuguese	29	7	36
Total	380	219	599

$$E_{ij} = \frac{n_{i.} \cdot n_{.j}}{n..}$$

$$E_{11} = \frac{380 \times 125}{599} = 79.3$$

$$E_{51} = \frac{380 \times 73}{599} = 46.3$$

$$E_{12} = \frac{219 \times 125}{599} = 45.7$$

$$E_{52} = \frac{219 \times 73}{599} = 26.69$$

$$E_{21} = \frac{380 \times 85}{599} = 53.92$$

$$E_{61} = \frac{380 \times 58}{599} = 36.71$$

$$E_{12} = \frac{85 \times 219}{599} = 31.08$$

$$E_{02} = \frac{219 \times 58}{599} = 21.11$$

$$E_{31} = \frac{380 \times 72}{599} = 45.68$$

$$E_{71} = \frac{380 \times 65}{599} = 23.71$$

$$E_{52} = \frac{219 \times 72}{599} = 26.32$$

$$E_{72} = \frac{219 \times 65}{599} = 23.71$$

$$E_{11} = \frac{380 \times 85}{599} = 53.92$$

$$E_{01} = \frac{380 \times 36}{599} = 22.84$$

$$E_{42} = \frac{219 \times 85}{599} = 31.08$$

$$E_{92} = \frac{219 \times 36}{599} = 13.16$$

$$\begin{aligned} X^2 &= \frac{(78 - 79.3)^2}{79.3} + \frac{(47 - 45.7)^2}{45.7} + \dots + \\ &\frac{(29 - 22.84)^2}{22.84} + \frac{(7 - 13.6)^2}{13.6} \\ &= 6.03 \end{aligned}$$

$$\chi^2_{(18-1)(2-1)}(1-0.05) = \chi^2_{(17)}(0.95) = 14.0671$$

$$P\text{-value} = P(X^2 > 6.03) = 0.5363$$

$\Rightarrow$  We don't reject  $H_0$

⑦ We test  $H_0$ : there is no relationship between major and grade

⇒ Chi-square test of homogeneity

$$\chi^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2_{((I-1)(J-1))}$$

Grade	Psychology	Biology	Other	Total
A	8	15	13	36
B	14	19	15	48
C	15	4	7	26
D-F	3	1	4	8
Total	40	39	39	118

$$E_{11} = \frac{40 \times 36}{118} = 12.20$$

$$E_{31} = \frac{40 \times 26}{118} = 8.81$$

$$E_{12} = \frac{39 \times 36}{118} = 11.90$$

$$E_{32} = \frac{39 \times 26}{118} = 8.59$$

$$E_{13} = \frac{39 \times 36}{118} = 11.90$$

$$E_{33} = \frac{39 \times 26}{118} = 8.59$$

$$E_{21} = \frac{40 \times 48}{118} = 16.27$$

$$E_{41} = \frac{40 \times 8}{118} = 2.64$$

$$E_{22} = \frac{39 \times 48}{118} = 15.96$$

$$E_{42} = \frac{39 \times 8}{118} = 2.64$$

$$E_{23} = \frac{39 \times 48}{118} = 15.96$$

$$E_{43} = \frac{39 \times 8}{118} = 2.64$$

$$\chi^2 = \frac{(8-12.2)^2}{12.2} + \frac{(15-11.9)^2}{11.9} + \dots + \frac{(4-2.64)^2}{2.64} = 12.18$$

$$\chi^2_{((4-1)(3-1))} (1-0.05) = \chi^2_{(6)} (0.95) = 12.5916$$

P-value = 0.0581

⇒ we don't reject  $H_0$

13

# children younger sister	# children older sister			Total
	0	1	2 >	
0	$n_{11}$	$n_{12}$	$n_{13}$	$n_{1.}$
1	$n_{21}$	$n_{22}$	$n_{23}$	$n_{2.}$
2 >	$n_{31}$	$n_{32}$	$n_{33}$	$n_{3.}$
Total	$n_{.1}$	$n_{.2}$	$n_{.3}$	$n_{..}$

a) The chi-square test of independence

$H_0: \pi_{ij} = \pi_{i.} \pi_{.j}$   $(i=1,2,3 \quad j=1,2,3)$  where  $\pi_{ij}$  are the multinomial cell probabilities of the counts  $n_{ij}$

Which is tested with  $\chi^2 = \sum_{i=1}^3 \sum_{j=1}^3 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = \sum_i \sum_j \frac{(n_{ij} - \frac{n_{i.} n_{.j}}{n_{..}})^2}{\frac{n_{i.} n_{.j}}{n_{..}}} \sim \chi^2_{(3-1)(3-1)}$

b)  $H_0: \sum \pi_{ij} = \sum \pi_{ji}$

$\Leftrightarrow \pi_{ij} = \pi_{ji} \quad i \neq j$

$$\Rightarrow \hat{\pi}_{12} = \hat{\pi}_{21} = \frac{n_{12} + n_{21}}{2n}$$

$$\hat{\pi}_{13} = \hat{\pi}_{31} = \frac{n_{13} + n_{31}}{2n}$$

$$\hat{\pi}_{23} = \hat{\pi}_{32} = \frac{n_{23} + n_{32}}{2n}$$

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

$$\chi^2 = \sum_i \sum_j \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = \sum_i \sum_j \frac{(n_{ij} - n \hat{\pi}_{ij})^2}{n \hat{\pi}_{ij}} = \sum_{i \neq j} \frac{(n_{ij} - \frac{n n_{ij}}{n})^2}{n \frac{n_{ij}}{n}} + \sum_{i \neq j} \frac{(n_{ij} - \frac{n(n_{ij} + n_{ji})}{2n})^2}{n \frac{(n_{ij} + n_{ji})}{2n}}$$

$$= \sum_{i \neq j} \frac{(n_{ij} - \frac{n_{ij} + n_{ji}}{2})^2}{\frac{n_{ij} + n_{ji}}{2}}$$

$$\chi^2 \sim \chi^2_{(2)}$$

df = number of free parameters under  $H_0$  - number of estimated parameter under  $H_1$   
 $= (9-1) - 6$   
 $= 2$

The null hypothesis are not equivalent, for example, if the younger sister had exactly the same number of children as the older a) would be false and b) true

(21)

	$\bar{D}$	$D$
$\bar{X}$	$\pi_{00}$	$\pi_{01}$
$X$	$\pi_{10}$	$\pi_{11}$

$$\text{odds}(D|X) = \frac{\pi_{11}}{\pi_{10}}$$

$$\text{odds}(D|\bar{X}) = \frac{\pi_{01}}{\pi_{00}}$$

$$\hat{\Delta} = \frac{n_{00}n_{11}}{n_{01}n_{10}}$$

$$\Delta = \frac{\text{odds}(D|X)}{\text{odds}(D|\bar{X})} = \frac{\pi_{11}\pi_{00}}{\pi_{01}\pi_{10}}$$

Normal      Diabetic

BB	49	39
Bb or bb	4	12

$$\hat{\Delta} = \frac{49 \cdot 12}{39 \cdot 4} = 3.77$$

The odds of contracting diabetes are increased by a factor of 3.77 if patient has Bb or bb allele

Chapter 14

②

X	0.34	1.38	-0.65	0.68	1.40	-0.88	-0.30	-1.18	0.50	-1.75
y	0.27	1.34	-0.53	0.35	1.28	-0.98	-0.71	-0.81	0.64	-1.59

a) Find  $y = a + bx$        $y = \beta_0 + \beta_1 x$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

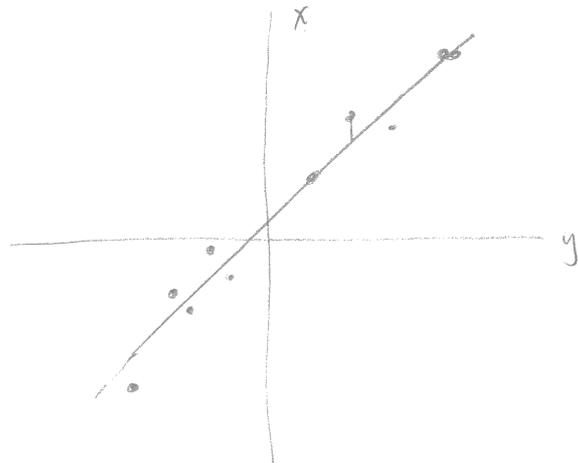
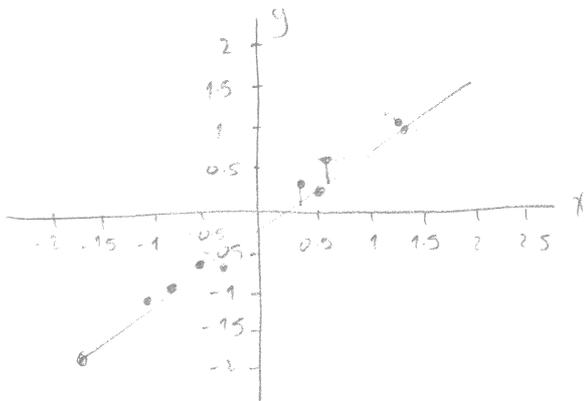
$$\bar{x} = -0.046$$

$$\bar{y} = -0.075$$

$$\hat{\beta}_1 = \frac{(0.34 + 0.046)(0.27 + 0.075) + \dots + (-1.75 + 0.046)(-1.59 + 0.075)}{(0.34 + 0.046)^2 + \dots + (-1.75 + 0.046)^2} = \frac{9.4177}{10.4130} = 0.9044$$

$$\hat{\beta}_0 = -0.075 - (0.9044)(-0.046) = -0.0334$$

$$\Rightarrow y = -0.0334 + 0.9044x$$



b) Find  $x = c + dy$

$$\hat{\beta}_0 = \bar{x} - \hat{\beta}_1 \bar{y} \quad \hat{\beta}_1 = \frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{\sum_i (y_i - \bar{y})^2}$$

$$\hat{\beta}_1 = \frac{(0.27 + 0.075)(0.34 + 0.046) + \dots + (-1.59 + 0.075)(-1.75 + 0.046)}{(0.27 + 0.075)^2 + \dots + (-1.59 + 0.075)^2} = \frac{9.4177}{8.9267} = 1.0550$$

$$\hat{\beta}_0 = -0.046 - (1.0550)(-0.075) = 0.0331$$

$$\Rightarrow x = 0.0331 + 1.0550y \Rightarrow c) \quad y = -0.0314 + 0.94179x \quad \text{so, not the same line as in a)}$$

c) because the distances (projections) to minimize in a) are not the same as in b) which one to use is not a problem, it's defined by which one is the independent (predictor) variable (in whose terms the properties, such as variance will be defined) and which one is the dependent (response) variable.

11) If  $\bar{x} = 0 \Rightarrow \text{cov}(\hat{\beta}_0, \hat{\beta}_1)$  under the assumptions of the standard statistical model.

Standard statistical model  $y_i = \beta_0 + \beta_1 x_i + e_i \quad i=1, \dots, n$   $E(e_i) = 0$   
 $\text{Var}(e_i) = \sigma^2$

under those assumptions  $\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} = \frac{-\sigma^2 n \bar{x}}{n \sum x_i - (\sum x_i)^2} = 0$   
 Th B. 1.9 5.9.9

$$\text{corr}(\hat{\beta}_0, \hat{\beta}_1) = \frac{\text{cov}(\hat{\beta}_0, \hat{\beta}_1)}{\sqrt{\text{Var}(\hat{\beta}_0) \text{Var}(\hat{\beta}_1)}} = 0$$

15) Find the least squares estimate of  $\beta$  for  $y = \beta x$  for  $(x_i, y_i) \quad i=1, 2, \dots, n$

The estimator will be found through  $\min S(\beta) = \sum_{i=1}^n (y_i - \beta x_i)^2$

$$\frac{\partial S}{\partial \beta} = -2 \sum_{i=1}^n x_i (y_i - \beta x_i) \quad \frac{\partial S}{\partial \beta} = 0 \Leftrightarrow \sum_{i=1}^n x_i (y_i - \beta x_i) = 0$$

$$\Leftrightarrow \sum x_i y_i = \beta \sum x_i^2 \Leftrightarrow \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

Another possibility is to find the mle (which coincides with  $\hat{\beta}$ ). <sup>\* solution in next page</sup>

19)  $\hat{Y} = X\hat{\beta} \quad \hat{e} = Y - X\hat{\beta} \quad \hat{\beta} = (X^T X)^{-1} X^T Y$

a) To show  $X^T \hat{e} = 0$

$$X^T (Y - X\hat{\beta}) = X^T Y - X^T X \hat{\beta} = X^T Y - X^T X (X^T X)^{-1} X^T Y = X^T Y - X^T Y = 0$$

$\Rightarrow$  each column of  $X$ ,  $x_i$  satisfy  $x_i^T \hat{e} = 0 \Rightarrow$  they are orthogonal

b) If the model contains an intercept term then  $X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1,p-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{n,p-1} \end{bmatrix}$

we know that  $X^T \hat{e} = 0$ , in particular  $0 = X_1^T \hat{e} = (1, 1, \dots, 1) \begin{pmatrix} \hat{e}_1 \\ \hat{e}_2 \\ \vdots \\ \hat{e}_n \end{pmatrix} = \sum_{i=1}^n \hat{e}_i$

$$\Rightarrow E\left(\sum_{i=1}^n \hat{e}_i\right) = E(0) = 0$$

$$* y_i - \beta x_i = \varepsilon_i \sim N(0, \sigma^2)$$

$$\Rightarrow y_i \sim (\beta x_i, \sigma^2)$$

$$L(\beta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \beta x_i)^2\right\}$$

$$l(\beta) = -\frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta x_i)^2$$

$$\frac{\partial l}{\partial \beta} = -\frac{1}{\sigma^2} \sum_{i=1}^n x_i (y_i - \beta x_i) = 0 \Leftrightarrow \sum_{i=1}^n x_i (y_i - \beta x_i) = 0$$

(35)

$$X = [X_1, X_2, X_1 + X_2]$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y = ([X_1, X_2, X_1 + X_2]^T [X_1, X_2, X_1 + X_2])^{-1} [X_1, X_2, X_1 + X_2]^T Y$$

$$= \begin{bmatrix} X_1^T \\ X_2^T \\ X_1^T + X_2^T \end{bmatrix} [X_1, X_2, X_1 + X_2]^{-1} \begin{bmatrix} X_1^T \\ X_2^T \\ X_1^T + X_2^T \end{bmatrix} Y$$

$$= \begin{bmatrix} X_1^T X_1 & X_1^T X_2 & X_1^T X_1 + X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 & X_2^T X_1 + X_2^T X_2 \\ X_1^T X_1 + X_2^T X_1 & X_1^T X_2 + X_2^T X_2 & X_1^T X_1 + X_1^T X_2 + X_2^T X_1 + X_2^T X_2 \end{bmatrix}^{-1} \begin{bmatrix} X_1^T \\ X_2^T \\ X_1^T + X_2^T \end{bmatrix} Y$$

$\Rightarrow (3) = (1) + (2) \Rightarrow X^T X$  doesn't have full rank  $\Rightarrow$  not possible to invert