

Suggested solutions for the 5th set of exercises

1. The conditional probability mass function is the product of the individual probability mass functions (by independence) divided by the probability mass function for the total count:

$$\begin{aligned}
 & P(N_1 = n_1, N_2 = n_2, \dots, N_c = n_c \mid \sum_{i=1}^c N_i = n) \\
 &= \frac{P(N_1 = n_1, N_2 = n_2, \dots, N_c = n_c)}{P(\sum_{i=1}^c N_i = n)} \\
 &= \frac{\prod_{i=1}^c e^{-\mu_i} \frac{\mu_i^{n_i}}{n_i!}}{e^{-\sum_{j=1}^c \mu_j} \frac{(\sum_{j=1}^c \mu_j)^n}{n!}} \\
 &= \frac{(\prod_{i=1}^c \frac{1}{n_i!}) \prod_{i=1}^c e^{-\mu_i} \mu_i^{n_i}}{\frac{1}{n!} e^{-\sum_{j=1}^c \mu_j} (\sum_{j=1}^c \mu_j)^{n_1} \dots (\sum_{j=1}^c \mu_j)^{n_c}} \\
 &= \frac{n!}{\prod_{i=1}^c n_i!} \prod_{i=1}^c e^{-\mu_i + \mu_i} \frac{\mu_i^{n_i}}{(\sum_{j=1}^c \mu_j)^{n_i}} \\
 &= \frac{n!}{\prod_{i=1}^c n_i!} \prod_{i=1}^c \left(\frac{\mu_i}{\sum_{j=1}^c \mu_j} \right)^{n_i} \\
 &= \frac{n!}{n_1! n_2! \dots n_c!} \pi_1^{n_1} \pi_2^{n_2} \dots \pi_c^{n_c},
 \end{aligned}$$

where $\pi_i = \mu_i / (\sum_{j=1}^c \mu_j)$. The required result was obtained.

2. Under the null hypothesis (proportions equal) the joint probability mass function is

$$\begin{aligned}
 & \binom{n_{1+}}{n_{11}} \pi^{n_{11}} (1 - \pi)^{n_{12}} \times \binom{n_{2+}}{n_{21}} \pi^{n_{21}} (1 - \pi)^{n_{22}} \\
 &= \binom{n_{1+}}{n_{11}} \binom{n_{2+}}{n_{21}} \pi^{n_{11} + n_{21}} (1 - \pi)^{n_{12} + n_{22}}
 \end{aligned}$$

by independence of the two samples. Under the null the marginal frequency n_{+1} likewise follows a Binomial distribution:

$$\binom{n}{n_{+1}} \pi^{n_{+1}} (1 - \pi)^{n - n_{+1}}.$$

Conditioning on the probability of the observed marginal frequency n_{+1} gives the hy-

pergeometric distribution asked for:

$$\frac{\binom{n_{1+}}{n_{11}} \binom{n_{2+}}{n_{21}} \pi^{n_{+1}} (1 - \pi)^{n_{+2}}}{\binom{n}{n_{+1}} \pi^{n_{+1}} (1 - \pi)^{n_{+2}}} = \frac{\binom{n_{1+}}{n_{11}} \binom{n_{2+}}{n_{21}}}{\binom{n}{n_{+1}}}.$$

If the column frequencies are not fixed, the same hypergeometric distribution arises, if the analysis is conditioned on the observed random column frequencies.

3.

a) The multivariate random variate of the exercise can be expressed as a c -vector:

$$\begin{bmatrix} N_1 \\ \vdots \\ N_c \end{bmatrix} = \sum_{i=1}^n \mathbf{Y}_i.$$

Above $\mathbf{Y}_i = [Y_{i1} \dots Y_{ic}]'$ is a multinomially distributed random vector from the i th (independent) experiment in a sequence of n experiments. The components of \mathbf{Y}_i are 0s except one which takes value 1 randomly according to the probabilities in the vector $\boldsymbol{\pi} = [\pi_1 \dots \pi_c]'$. Then $\sum_{j=1}^c Y_{ij} = 1$, $Y_{ij} Y_{ik} = 0$, if $j \neq k$,

$$E(Y_{ij}) = \pi_j \times 1 + (1 - \pi_j) \times 0 = \pi_j = E(Y_{ij}^2)$$

and

$$E(Y_{ij} Y_{ik}) = 0, \quad \text{if } j \neq k.$$

Thus

$$E(\mathbf{Y}_i) = \boldsymbol{\pi}$$

and

$$\text{cov}(\mathbf{Y}_i) = \boldsymbol{\Sigma}.$$

The jk th component of the covariance matrix above is $\{\boldsymbol{\Sigma}\}_{jk} = \sigma_{jk}$, where

$$\sigma_{jj} = \text{var}(Y_{ij}) = E(Y_{ij}^2) - [E(Y_{ij})]^2 = \pi_j - \pi_j^2 = \pi_j(1 - \pi_j)$$

($j = k$) and

$$\sigma_{jk} = \text{cov}(Y_{ij}, Y_{ik}) = E(Y_{ij} Y_{ik}) - E(Y_{ij}) E(Y_{ik}) = 0 - \pi_j \pi_k = -\pi_j \pi_k$$

($j \neq k$). The covariance matrix of vector $\sum_{i=1}^n \mathbf{Y}_i$ is

$$\text{cov}\left(\sum_{i=1}^n \mathbf{Y}_i\right) = \sum_{i=1}^n \text{cov}(\mathbf{Y}_i) = n\boldsymbol{\Sigma}$$

by independence. Hence the covariance between the frequencies in categories j and k is

$$\text{cov}(N_j, N_k) = \{n\Sigma\}_{jk} = n\sigma_{jk} = -n\pi_j\pi_k.$$

b) Because of point a) it is the case that

$$\text{cov}(N_j, N_k) = -n\pi_j\pi_k.$$

From Exercise 4.1:

$$\text{var}(N_j) = n\pi_j(1 - \pi_j).$$

Substituting these to the definition of correlation gives

$$\begin{aligned} \text{cor}(N_j, N_k) &= \frac{\text{cov}(N_j, N_k)}{\sqrt{\text{var}(N_j)}\sqrt{\text{var}(N_k)}} \\ &= \frac{-\pi_j\pi_k}{\sqrt{\pi_j(1 - \pi_j)}\sqrt{\pi_k(1 - \pi_k)}} \\ &= \frac{-\pi_j\pi_k}{\sqrt{\pi_j(1 - \pi_j)\pi_k(1 - \pi_k)}}. \end{aligned}$$

c) Here $c = 2$, $\pi_2 = 1 - \pi_1$ and $N_2 = n - N_1$. The correlation between frequencies N_1 and N_2 is now -1 :

$$\begin{aligned} \text{cor}(N_1, N_2) &= \frac{-\pi_1\pi_2}{\sqrt{\pi_1(1 - \pi_1)\pi_2(1 - \pi_2)}} \\ &\stackrel{\pi_2=1-\pi_1}{=} \frac{-\pi_1(1 - \pi_1)}{\sqrt{\pi_1(1 - \pi_1)(1 - \pi_1)\pi_1}} \\ &= \frac{-\pi_1(1 - \pi_1)}{\sqrt{[\pi_1(1 - \pi_1)]^2}} \\ &= -1. \end{aligned}$$

The intuition is evident: There is a perfect linear relation between the frequencies $N_2 = n - N_1$. If the other increases, the other decreases, and *vice versa*.

4.

a) The observations are binomially distributed in the two samples. The probability mass functions of the first and second sample are

$$\binom{n_{1+}}{n_{11}} \pi_1^{n_{11}} (1 - \pi_1)^{n_{1+} - n_{11}} = \binom{n_{1+}}{n_{11}} \pi_1^{n_{11}} (1 - \pi_1)^{n_{12}}$$

and

$$\binom{n_{2+}}{n_{21}} \pi_2^{n_{21}} (1 - \pi_2)^{n_{2+} - n_{21}} = \binom{n_{2+}}{n_{21}} \pi_2^{n_{21}} (1 - \pi_2)^{n_{22}},$$

respectively. Due to the independence of the samples, the joint probability mass function is

$$\binom{n_{1+}}{n_{11}} \pi_1^{n_{11}} (1 - \pi_1)^{n_{12}} \times \binom{n_{2+}}{n_{21}} \pi_2^{n_{21}} (1 - \pi_2)^{n_{22}}.$$

Taking the log of it and deleting additive terms, which do not involve π_1 or π_2 , yields the log-likelihood function

$$l(\boldsymbol{\pi}) = n_{11} \log \pi_1 + n_{12} \log(1 - \pi_1) + n_{21} \log \pi_2 + n_{22} \log(1 - \pi_2).$$

b) The partial derivatives of the log-likelihood function with respect to π_1 and π_2 are

$$\frac{\partial}{\partial \pi_1} l(\boldsymbol{\pi}) = n_{11} \frac{\partial}{\partial \pi_1} \log \pi_1 + n_{12} \frac{\partial}{\partial \pi_1} \log(1 - \pi_1) = \frac{n_{11}}{\pi_1} - \frac{n_{12}}{1 - \pi_1}$$

and

$$\frac{\partial}{\partial \pi_2} l(\boldsymbol{\pi}) = n_{21} \frac{\partial}{\partial \pi_2} \log \pi_2 + n_{22} \frac{\partial}{\partial \pi_2} \log(1 - \pi_2) = \frac{n_{21}}{\pi_2} - \frac{n_{22}}{1 - \pi_2},$$

respectively. It follows that

$$\nabla l(\boldsymbol{\pi}) = \left[\frac{n_{11}}{\pi_1} - \frac{n_{12}}{1 - \pi_1} \quad \frac{n_{21}}{\pi_2} - \frac{n_{22}}{1 - \pi_2} \right]'$$

c) Under the null hypothesis of $\pi_1 = \pi_2$ the restricted MLE for both parameters is $\hat{\pi}_{1,0} = \hat{\pi}_{2,0} = n_{+1}/n$. The partial derivatives evaluated at this restricted MLE are

$$\begin{aligned} \nabla l(\hat{\boldsymbol{\pi}}_0) &= \left[\frac{n_{11}}{n_{+1}/n} - \frac{n_{12}}{1 - n_{+1}/n} \quad \frac{n_{21}}{n_{+1}/n} - \frac{n_{22}}{1 - n_{+1}/n} \right]' \\ &= \left[n \left(\frac{n_{11}}{n_{+1}} - \frac{n_{12}}{n_{+2}} \right) \quad n \left(\frac{n_{21}}{n_{+1}} - \frac{n_{22}}{n_{+2}} \right) \right]' \\ &= \left[n \frac{n_{11}n_{+2} - n_{12}n_{+1}}{n_{+1}n_{+2}} \quad n \frac{n_{21}n_{+2} - n_{22}n_{+1}}{n_{+1}n_{+2}} \right]' \\ &= \left[n \frac{n_{11}(n_{12} + n_{22}) - n_{12}(n_{11} + n_{21})}{n_{+1}n_{+2}} \quad n \frac{n_{21}(n_{12} + n_{22}) - n_{22}(n_{11} + n_{21})}{n_{+1}n_{+2}} \right]' \\ &= \left[n \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{+1}n_{+2}} \quad n \frac{n_{12}n_{21} - n_{11}n_{22}}{n_{+1}n_{+2}} \right]'. \end{aligned}$$

d)

$$\begin{aligned} \frac{\partial^2}{\partial \pi_1^2} l(\boldsymbol{\pi}) &= \frac{-n_{11}}{\pi_1^2} - \frac{-(-1)n_{12}}{(1 - \pi_1)^2} = - \left[\frac{n_{11}}{\pi_1^2} + \frac{n_{12}}{(1 - \pi_1)^2} \right], \\ \frac{\partial^2}{\partial \pi_2^2} l(\boldsymbol{\pi}) &= \frac{-n_{21}}{\pi_2^2} - \frac{-(-1)n_{22}}{(1 - \pi_2)^2} = - \left[\frac{n_{21}}{\pi_2^2} + \frac{n_{22}}{(1 - \pi_2)^2} \right], \end{aligned}$$

and

$$\frac{\partial^2}{\partial \pi_1 \partial \pi_2} l(\boldsymbol{\pi}) = \frac{\partial^2}{\partial \pi_2 \partial \pi_1} l(\boldsymbol{\pi}) = 0.$$

Hence

$$- \begin{bmatrix} \frac{\partial^2}{\partial \pi_1^2} l(\boldsymbol{\pi}) & \frac{\partial^2}{\partial \pi_1 \partial \pi_2} l(\boldsymbol{\pi}) \\ \frac{\partial^2}{\partial \pi_2 \partial \pi_1} l(\boldsymbol{\pi}) & \frac{\partial^2}{\partial \pi_2^2} l(\boldsymbol{\pi}) \end{bmatrix} = \begin{bmatrix} \frac{n_{11}}{\pi_1^2} + \frac{n_{12}}{(1-\pi_1)^2} & 0 \\ 0 & \frac{n_{21}}{\pi_2^2} + \frac{n_{22}}{(1-\pi_2)^2} \end{bmatrix}.$$

e) The Fisher information matrix is

$$\begin{aligned} \mathbf{I}(\boldsymbol{\pi}) &= \mathbb{E} \begin{bmatrix} \frac{n_{11}}{\pi_1^2} + \frac{n_{12}}{(1-\pi_1)^2} & 0 \\ 0 & \frac{n_{21}}{\pi_2^2} + \frac{n_{22}}{(1-\pi_2)^2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n_{1+}\pi_1}{\pi_1^2} + \frac{n_{1+}(1-\pi_1)}{(1-\pi_1)^2} & 0 \\ 0 & \frac{n_{2+}\pi_2}{\pi_2^2} + \frac{n_{2+}(1-\pi_2)}{(1-\pi_2)^2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n_{1+}}{\pi_1} + \frac{n_{1+}}{1-\pi_1} & 0 \\ 0 & \frac{n_{2+}}{\pi_2} + \frac{n_{2+}}{1-\pi_2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n_{1+}(1-\pi_1) + n_{1+}\pi_1}{\pi_1(1-\pi_1)} & 0 \\ 0 & \frac{n_{2+}(1-\pi_2) + n_{2+}\pi_2}{\pi_2(1-\pi_2)} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n_{1+}}{\pi_1(1-\pi_1)} & 0 \\ 0 & \frac{n_{2+}}{\pi_2(1-\pi_2)} \end{bmatrix}. \end{aligned}$$

The second equality is due to $n\pi$ being the expected value of the Binomial distribution $\text{Bin}(n, \pi)$.

The inverse matrix of the Fisher information matrix is

$$\mathbf{I}(\boldsymbol{\pi})^{-1} = \begin{bmatrix} \frac{\pi_1(1-\pi_1)}{n_{1+}} & 0 \\ 0 & \frac{\pi_2(1-\pi_2)}{n_{2+}} \end{bmatrix}.$$

Evaluating it at the restricted MLE $\hat{\pi}_0$ yields

$$\begin{aligned} \mathbf{I}(\hat{\pi}_0)^{-1} &= \begin{bmatrix} \frac{n+1}{n} \left(1 - \frac{n+1}{n}\right) & 0 \\ n_{1+} & \frac{n+1}{n} \left(1 - \frac{n+1}{n}\right) \\ 0 & n_{2+} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n+1n_{+2}}{n^2n_{1+}} & 0 \\ 0 & \frac{n+1n_{+2}}{n^2n_{2+}} \end{bmatrix}. \end{aligned}$$

f) On the grounds of the previous calculations

$$\begin{aligned} \nabla l(\hat{\pi}_0)' \mathbf{I}(\hat{\pi}_0)^{-1} \nabla l(\hat{\pi}_0) &\stackrel{\text{c) and e)}}{=} \begin{bmatrix} n \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{+1}n_{+2}} & n \frac{n_{12}n_{21} - n_{11}n_{22}}{n_{+1}n_{+2}} \end{bmatrix} \times \\ &\begin{bmatrix} \frac{n+1n_{+2}}{n^2n_{1+}} & 0 \\ 0 & \frac{n+1n_{+2}}{n^2n_{2+}} \end{bmatrix} \begin{bmatrix} n \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{+1}n_{+2}} \\ n \frac{n_{12}n_{21} - n_{11}n_{22}}{n_{+1}n_{+2}} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n_{11}n_{22} - n_{12}n_{21}}{nn_{1+}} & \frac{n_{12}n_{21} - n_{11}n_{22}}{nn_{2+}} \end{bmatrix} \times \\ &\begin{bmatrix} n \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{+1}n_{+2}} \\ n \frac{n_{12}n_{21} - n_{11}n_{22}}{n_{+1}n_{+2}} \end{bmatrix} \\ &= \frac{(n_{11}n_{22} - n_{12}n_{21})^2}{n_{1+}n_{+1}n_{+2}} + \frac{(n_{11}n_{22} - n_{12}n_{21})^2}{n_{2+}n_{+1}n_{+2}} \\ &= \frac{(n_{11}n_{22} - n_{12}n_{21})^2 (n_{1+} + n_{2+})}{n_{1+}n_{2+}n_{+1}n_{+2}} \\ &= \frac{n(n_{11}n_{22} - n_{12}n_{21})^2}{n_{1+}n_{2+}n_{+1}n_{+2}}. \end{aligned}$$

The formula for the statistic for testing independence derived in Exercise 3.6 was obtained. It equals test statistic z_s^2 according to Exercise 4.5. Test statistics z_s^2 and X^2 are score statistics.

Extra comments:

- Cox and Hinkley (1978, 133–134) give an alternative proof.¹ They additionally prove that the Wald statistic is approximately

$$\sum_{i=1}^2 \sum_{j=1}^2 \frac{(n_{ij} - \hat{\mu}_{ij})^2}{n_{ij}}$$

¹D.R. Cox ja D. Hinkley (1978): *Problems and Solutions in Theoretical Statistics*. Chapman and Hall, Lontoo.

in the present circumstance.

- The likelihood ratio statistic could be calculated as in Exercise 4.4 (the likelihood ratio statistic is the same despite that the method of sampling is different).
- The outcome of the exercise can be generalised:
 1. The score statistic and Pearson's test statistic for independence X^2 or

$$\sum_{i=1}^J \sum_{j=1}^K \frac{(n_{ij} - \hat{\mu}_{ij})^2}{\hat{\mu}_{ij}}$$

equal also in the case of a $J \times K$ ($J \geq 2$ and $K \geq 2$) contingency table (*op. cit.*). Above $\hat{\mu}_{ij} = n_{i+}n_{+j}/n$ (in obvious notation). The Wald statistic generalises correspondingly (*op. cit.*).

2. According to Silvey (1975, 120)² the score statistic and X^2

$$\sum_{i=1}^c \frac{(n_i - \hat{\mu}_i)^2}{\hat{\mu}_i}$$

match in general in the context of testing of the parameters of a c dimensional multinomial distribution. Here $\hat{\mu}_i$ is estimated under the restrictions which apply under the null hypothesis. Cox and Hinkley (1974, 315–316 and 325–327)³ point out related results.

3. Agresti (2013, 22)⁴ proves that in the case of a multinomial distribution with c -categories

$$n(\hat{\pi} - \pi_0)' \Sigma_0^{-1} (\hat{\pi} - \pi_0) = \sum_{i=1}^c \frac{(n_i - \mu_i)^2}{\mu_i}.$$

Here $\hat{\pi} = [\hat{\pi}_1 \dots \hat{\pi}_{c-1}]'$ is the (unrestricted) MLE for the parameter vector π which is π_0 under the null hypothesis, Σ_0 is the asymptotic covariance matrix of $\sqrt{n}\hat{\pi}$ under the null hypothesis and μ_i is the expected frequency of the i th cell under the null hypothesis.

4. Agresti (*op. cit.*, 78–80) discusses construction of score confidence intervals for the differences of two proportions and points out R code for it (http://www.stat.ufl.edu/~aa/cda/R_web.pdf; read 8.10.2015).

²S.D. Silvey (1975): *Statistical Inference*. Chapman and Hall. London.

³D.R. Cox and D.V. Hinkley (1974): *Theoretical Statistics*. Chapman and Hall. London.

⁴A. Agresti (2013): *Categorical Data Analysis, 3rd edition*. Wiley. Hoboken, NJ.

5.

a) If marginal homogeneity applies then $\pi_{1+} = \pi_{+1}$. It is then also the case that $\pi_{12} = \pi_{21}$:

$$\begin{aligned} 0 &= \pi_{1+} - \pi_{+1} \\ &= \pi_{11} + \pi_{12} - (\pi_{11} + \pi_{21}) \\ &= \pi_{12} - \pi_{21}. \end{aligned}$$

b) If $\pi_{1+} = \pi_{+1}$ then also $\pi_{2+} = \pi_{+2}$:

$$\begin{aligned} \pi_{2+} - \pi_{+2} &= \pi_{21} + \pi_{22} - (\pi_{12} + \pi_{22}) \\ &= \pi_{21} - \pi_{12} \\ &\stackrel{a)}{=} -(\pi_{1+} - \pi_{+1}) \\ &= 0. \end{aligned}$$

The last equality follows from the assumption of marginal homogeneity.

c) Marginal homogeneity ($\pi_{1+} = \pi_{+1}$ ja $\pi_{2+} = \pi_{+2}$) does not imply $\pi_{11} = \pi_{22}$. The counter examples below make the point.

		Y		
		y_1	y_2	Σ
X	x_1	0	0,1	0,1
	x_2	0,1	0,8	0,9
Σ		0,1	0,9	1

		Y		
		y_1	y_2	Σ
X	x_1	0	0,1	0,1
	x_2	0,1	0,8	0,9
Σ		0,1	0,9	1

		Y		
		y_1	y_2	Σ
X	x_1	0,3	0,15	0,45
	x_2	0,15	0,4	0,55
Σ		0,45	0,55	1

		Y		
		y_1	y_2	Σ
X	x_1	0,30	0,25	0,55
	x_2	0,25	0,20	0,45
Σ		0,55	0,45	1

		Y		
		y_1	y_2	Σ
X	x_1	0,15	0,40	0,55
	x_2	0,40	0,05	0,45
Σ		0,55	0,45	1

6. The alternative formulations of McNemar's test statistic can be derived easily:

$$\begin{aligned}
\frac{n_{12} - 0,5 \times n^*}{\sqrt{n^* \times 0,5 \times 0,5}} &= \frac{n_{12} - 0,5 \times (n_{12} + n_{21})}{\sqrt{(n_{12} + n_{21}) \times 0,5 \times 0,5}} \\
&= \frac{0,5(n_{12} - n_{21})}{0,5\sqrt{n_{12} + n_{21}}} \\
&= \frac{n_{12} - n_{21}}{\sqrt{n_{12} + n_{21}}} \\
&= \frac{0,5(n_{12} - n_{21})}{0,5\sqrt{n_{12} + n_{21}}} \\
&= \frac{-[n_{21} - 0,5(n_{12} + n_{21})]}{\sqrt{(n_{12} + n_{21}) \times 0,5 \times 0,5}} \\
&= \frac{-(n_{21} - 0,5n^*)}{\sqrt{n^* \times 0,5 \times 0,5}}.
\end{aligned}$$

McNemar's test statistic can be calculated by means of the frequency n_{12} or n_{21} .