

# Assignment 4

Wednesday, 26 October 2022 7:51 PM

Q2

a. We have a multinomial distribution with

$$K=4, n = X_1 + X_2 + X_3 + X_4,$$

$$p_1, p_2, p_3, \text{ and } p_4 = 1 - p_1 - p_2 - p_3$$

The likelihood function of the multinomial function is given by

$$L(\boldsymbol{p} | \underline{X}) = \binom{n}{\underline{X}} p_1^{x_1} \times \dots \times p_K^{x_K}$$

So, for our  $K=4$ , we have

$$L(\boldsymbol{p} | \underline{X}) = \binom{n}{\underline{X}} p_1^{x_1} \cdot p_2^{x_2} \cdot p_3^{x_3} \cdot (1 - p_1 - p_2 - p_3)^{x_4}$$

or

$$L(\boldsymbol{p} | \underline{X}) = \frac{n!}{x_1! x_2! x_3! x_4!} p_1^{x_1} \cdot p_2^{x_2} \cdot p_3^{x_3} \cdot (1 - p_1 - p_2 - p_3)^{x_4}$$

b. Under the null hypothesis, we have

$$H_0: p_2 = p_3$$

Thus,  $p_1$  is a free parameter

$p_2/p_3$  count as one free parameter

and  $p_4 = 1 - p_1 - p_2 - p_3 = 1 - p_1 - 2p_2$  is not a free parameter.

i.e. There are two free parameters under the null.

c. Under the alternative

$$H_1: p_2 \neq p_3$$

we will have  $p_1, p_2$ , and  $p_3$  as free parameters

So, we have 3 free parameters under the alternative.

d. Under the null, the likelihood becomes

$$L(\boldsymbol{p} | \underline{X}) = n! \frac{p_1^{x_1}}{x_1!} \frac{p_2^{x_2+x_3}}{x_2! x_3!} \frac{p_4^{x_4}}{x_4!}$$

And the log likelihood is

$$\log L(\boldsymbol{p}) = \log(n!) + x_1 \log p_1 + (x_2 + x_3) \log p_2 + x_4 \log p_4 - \log(x_1! x_2! x_3! x_4!)$$

We note that  $p_1 + p_2 + p_3 + p_4 = 1$ , and  $n = X_1 + X_2 + X_3 + X_4$

Introduce a Lagrange constraint,  $\lambda$ , such that

$$\sum_{i=1}^4 p_i = \sum_{i=1}^4 \frac{x_i}{\lambda}$$

$$\begin{aligned} \lambda &= \frac{1}{\lambda} \sum_{i=1}^4 x_i \\ &= \frac{n}{\lambda} \Rightarrow \lambda = n \end{aligned}$$

Putting this as a constraint to the log likelihood, we get that

Putting this as a constraint to the log likelihood,  
we get that

$$\log L(\boldsymbol{\hat{p}}, \lambda) = \log L(\boldsymbol{\hat{p}} | \mathbf{x}) + \lambda \left( 1 - \sum_{i=1}^4 p_i \right)$$

Next, we want to find the MLE (i.e.  $\frac{\partial \log L}{\partial \lambda} = 0$ )

for  $p_1$ :

$$\frac{\partial \log L(\boldsymbol{\hat{p}}, \lambda)}{\partial p_1} = \frac{x_1}{\hat{p}_1} - \lambda = 0$$

$$\Rightarrow \frac{x_1}{\hat{p}_1} = \lambda = n$$

$$\Rightarrow \hat{p}_1 = \frac{x_1}{n}$$

An identical result can be obtained for each  $(\hat{p}_i, x_i)$  pair. However,  $p_2 = p_3 \Rightarrow \hat{p}_2 = \hat{p}_3$   
under the null, and so

$$\hat{p}_2 + \hat{p}_3 = \frac{x_2}{n} + \frac{x_3}{n}$$

$$2\hat{p}_2 = 2\hat{p}_3 = \frac{x_2 + x_3}{n}$$

$$\Rightarrow \hat{p}_2 = \hat{p}_3 = \frac{x_2 + x_3}{2n}$$

c. Under the null hypothesis, we'd expect that,

$$\hat{p}_i = \frac{x_i}{n} \Rightarrow x_i = \hat{p}_i \cdot n$$

for each cell. That is

$$x_1 = \hat{p}_1 \cdot n; x_2 = x_3 = \hat{p}_2 \cdot n = \hat{p}_3 \cdot n; x_4 = \hat{p}_4 \cdot n$$

$$(or \quad x_4 = (1 - \hat{p}_1 - 2\hat{p}_2) \cdot n)$$

f. We have McNemar's test as

$$\frac{(x_2 - x_3)^2}{x_2 + x_3}$$

and Pearson's  $\chi^2$  test as

$$\begin{aligned} \chi^2 &= \sum_{\text{cells}} \frac{(O_i - E_i)^2}{E_i} \\ &= \frac{(x_1 - \hat{p}_1 \cdot n)^2}{\frac{x_1}{n} \cdot n} + \frac{(x_2 - \hat{p}_2 \cdot n)^2}{\frac{x_2}{n} \cdot n} + \frac{(x_3 - \hat{p}_3 \cdot n)^2}{\frac{x_3}{n} \cdot n} + \frac{(x_4 - \hat{p}_4 \cdot n)^2}{\frac{x_4}{n} \cdot n} \\ &= \frac{(x_1 - \frac{x_1}{n} \cdot n)^2}{n} + \frac{(x_2 - \frac{x_2 + x_3}{2n} \cdot n)^2}{\frac{x_2 + x_3}{2n} \cdot n} + \frac{(x_3 - \frac{x_2 + x_3}{2n} \cdot n)^2}{\frac{x_2 + x_3}{2n} \cdot n} + \frac{(x_4 - \frac{x_4}{n} \cdot n)^2}{\frac{x_4}{n} \cdot n} \\ &= \frac{(x_2 - \frac{x_2}{2} - \frac{x_3}{2})^2}{\frac{x_2}{2} \cdot \frac{x_3}{2}} + \frac{(x_3 - \frac{x_2}{2} - \frac{x_3}{2})^2}{\frac{x_2}{2} \cdot \frac{x_3}{2}} \\ &= \frac{(x_2 - x_3)^2}{4x_2} + \frac{(x_3 - x_2)^2}{4x_3} \end{aligned}$$

Now, notice that

$$x_2 = \hat{p}_2 \cdot n = \frac{x_2 + x_3}{2} = \hat{p}_3 \cdot n = x_3$$

$$\Rightarrow \chi^2 = \frac{(x_2 - x_3)^2 + (x_3 - x_2)^2}{4 \cdot \frac{(x_2 + x_3)}{2}}$$

$$= \frac{x_2^2 - 2x_2x_3 + x_3^2 + x_3^2 - 2x_2x_3 + x_2^2}{(\frac{x_2 + x_3}{2})}$$

$$= \frac{x_2^2 - 2x_2x_3 + x_3^2}{x_2 + x_3}$$

$$\begin{aligned}
 &= \frac{x_2^2 - 2x_2x_3 + x_3^2}{x_2+x_3} \\
 &= \frac{(x_2 - x_3)^2}{x_2+x_3}
 \end{aligned}$$

g. Since under the null,  $x_2$  and  $x_3$  are described by the same MLE, they should have the same value and so there are only two free parameters,  $p_1$  and  $p_2$  (or  $p_3$ ). Hence,  $dof = n_{\text{free}} - 1 = 2 - 1 = 1$ , where  $n_{\text{free}}$  is the number of free parameters.

h. Since we're dealing with  $dof = 1$ , and the tabulated data gives

$$\begin{aligned}
 \chi^2 &= \frac{(x_2 - x_3)^2}{x_2+x_3} \\
 &= \frac{(15-4)^2}{15+4} \approx 6.368
 \end{aligned}$$

our p-value is given by

$$P = P(\chi^2 \geq 6.368)$$

In python, we can calculate this by

`print(1 - scipy.stats.chi2.cdf(6.368, 1))`

which yields

$$p \approx 0.0116$$

This is strong evidence that  $p_2 \neq p_3$ , and that the null hypothesis doesn't explain the phenomena i.e. that  $n_{\text{willing-willing}} \neq n_{\text{unwilling-willing}}$ .

i. Since one cell had  $< 5$  counts, we need to use Fisher's exact test.

Firstly, we can consider the row and column totals fixed under the null, as we essentially have only one degree of freedom (with two free parameters).

As such, we can generate multiple tables as any should be equally likely under the null.

We used the python code below to analyse this

```

23 #q2i
24 row1 = 52; row2 = 30
25 col1 = 41; col2 = 41
26 iters = 10**6
27 counts = 0
28
29 for i in range(iters):
30     x1 = np.random.randint(0, min(row1 + 1, col1))
31     x2 = row1 - x1
32     x3 = col1 - x1
33     x4 = col2 - x2
34     if (x2 - x3)**2 / (x2 + x3) >= 6.368:
35         counts += 1
36
37 print(f"p-value from {iters} trials is ", counts / iters)

```

This produced a p-value of

$p \approx 0.098$  (on average)  
from  $10^6$  table variations.

$p \approx 0.098$  (on average)  
from  $10^6$  table variations.  
This says that there is no significant  
evidence against the null hypothesis via  
Fisher's exact test. (i.e. our conclusion from  
h) was likely wrong).

Q3

We have the numerical data,

maths terms	incorrect solutions	correct solutions	total
$\mu$	298	165	463
$X$	146	264	410
$S$ or $S^2$	119	16	135
$\leq$ or $\geq$ or $<$ or $>$	106	328	434
1.96 or $-1.96$	29	203	232
$\sqrt{n}$	133	165	298
$\infty$	29	30	59
Total	860	1171	2031

a. In this situation, our null and alternative hypotheses are

$H_0$ : The relative frequency of each math term is the same across correct/incorrect solutions  
 $H_1$ : The relative frequency of at least one math term differs between correct/incorrect solutions.

We are interested here in the counts of each term (for correct and incorrect solutions), as well as each row/column total.

Notice that we'll have  $(n_{\text{rows}} - 1) \times (n_{\text{cols}} - 1)$  degrees of freedom. That is,

$$dof = (7-1) \times (2-1) = 6$$

So we can calculate an approximate p-value via

$$\text{p-value} = P(\chi^2_6 \geq \chi^2)$$

where  $\chi^2$  is Pearson's  $\chi^2$  statistic, given by

$$\chi^2 = \sum_{\text{cells}} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where  $O_{ij}$  is the observed value in the  $(i, j)^{\text{th}}$  cell,  
and

$$E_{ij} = \frac{(\sum O_{i\cdot}) \times (\sum O_{\cdot j})}{N}$$

where  $N$  is the sum of the cell counts.

Using the following python code,

```
col1 = np.array([298, 146, 119, 106, 29, 133, 29])
col2 = np.array([165, 264, 16, 328, 203, 165, 30])
total = sum(col1) + sum(col2)
chi2 = 0
for i in range(len(col1)):
    for col in [col1, col2]:
        Eij = (sum(col) * sum([col1[i], col2[i]])) / total
        chi2 += (col[i] - Eij)**2 / Eij
pval = 1 - stats.chi2.cdf(chi2, 6)
print(f"We get a value of {chi2 = }, and a corresponding {pval = }")
```

we obtain a Pearson's  $\chi^2$  of

$$\chi^2 \approx 359.06$$

and a p-value  $< 1e-16$  (almost certain significance)  
(python displays a pval=0.0, so the p-value is

$$\approx 0.01.06$$

and a p-value  $< 1e-16$  (almost certain significance)  
(Python displays a pval = 0.0, so the p-value is below floating point precision).

So, we have incredibly strong evidence that the usage of at least one mathematical term differs between correct and incorrect solutions.

b. For the 8 solutions that got 0/4 marks, we can employ a similar process to part a)

maths terms	incorrect solution								96 50 35 34 8 39 8
	1	2	3	4	5	6	7	8	
$\mu$	12	11	12	12	13	11	13	12	96
$\bar{X}$	6	6	7	6	5	8	7	5	50
$S$ or $S^2$	5	5	5	4	5	3	4	4	35
$\leq$ or $\geq$ or $<$ or $>$	4	4	6	2	4	2	6	6	34
1.96 or -1.96	2	0	0	0	2	0	2	2	8
$\sqrt{n}$	5	6	4	5	4	4	6	5	39
$\infty$	2	0	2	2	0	0	0	2	8
	36	32	36	31	33	28	38	36	270

In this case, our null and alternative hypotheses will be

$H_0$ : The relative frequency of math terms is consistent across the 8 answers

$H_1$ : There is at least one solution that uses a different relative frequency of math terms.

We can use Pearson's  $\chi^2$  test again, with

$$df = (\text{rows} - 1) \times (\text{cols} - 1) = (7-1) \times (8-1) = 42$$

and  $p\text{-value} = P(\chi^2_{42} \geq \chi^2)$

Using the following python code,

```
col1 = np.array([12, 6, 5, 4, 2, 5, 2])
col2 = np.array([11, 6, 5, 4, 0, 6, 0])
col3 = np.array([12, 7, 5, 6, 0, 4, 2])
col4 = np.array([12, 6, 4, 2, 0, 5, 2])
col5 = np.array([13, 5, 5, 4, 2, 4, 0])
col6 = np.array([11, 8, 3, 2, 0, 4, 0])
col7 = np.array([13, 7, 4, 6, 2, 6, 0])
col8 = np.array([12, 5, 4, 6, 2, 5, 2])
total = np.sum([[col1, col2, col3, col4, col5, col6, col7, col8]])
chi2 = 0
for i in range(len(col1)):
    for col in [col1, col2, col3, col4, col5, col6, col7, col8]:
        Eij = (sum(col) * sum([col1[i], col2[i], col3[i], col4[i], col5[i], col6[i], col7[i], col8[i]])) / total
        chi2 += (col[i] - Eij)**2 / Eij
pval = 1 - stats.chi2.cdf(chi2, 42)
print(f"We get a value of {chi2 = }, and a corresponding {pval = }")
```

We obtain  $\chi^2 \approx 22.08$  and  $p\text{-val} \approx 0.995$

This is incredibly weak evidence against the null hypothesis, and so we conclude that all 8 incorrect solutions use the same relative frequency of math terms.

c. The p-value in part b) is definitely suspicious, as it is very large, and we would not expect that from independent answers.

The cause that first comes to mind (in the context of an assignment solution) is plagiarism amongst at least three students (and possibly more when looking at other incorrect solutions).

d. Since so many of the cells in the table have less than 5 counts, we don't meet the required conditions to definitively use Pearson's  $\chi^2$  test. As such, we would need to perform

have less than 5 counts, we don't meet the required conditions to definitively use Pearson's  $\chi^2$  test. As such, we would need to perform Fisher's exact test on the observed data.

Q1

- a. Using a short python script with numpy, the data can be summarised by

$$\mu \approx 3.096 \quad \sigma \approx 1.59 \\ \text{median} = 3 \quad \text{mode} = 5 \quad (52 \text{ counts})$$

- b. For our hypothesis test, we propose

$H_0$ : the proportion of each opinion is consistent between regular people and influencers

$H_1$ : there is at least one difference in opinion proportions between regular people and influencers.

To test these hypotheses, we can employ Pearson's  $\chi^2$  statistic. Using the following python script,

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats

# q1a
data = pd.read_csv("likes.csv", delimiter=',')
print('mean = ', np.mean(data['rating']), ' sd = ', np.std(data['rating']),
      ' median = ', np.median(data['rating']),
      ' mode = ', stats.mode(data['rating']))

# q1b
InfAnswers = data[data['user'] == 'influencer']
RegAnswers = data[data['user'] == 'regular']
InfCounts, RegCounts = np.zeros(5), np.zeros(5)
for i in range(len(InfAnswers['rating'])):
    InfCounts[InfAnswers['rating'].to_numpy()[i] - 1] += 1
for i in range(len(RegAnswers['rating'])):
    RegCounts[RegAnswers['rating'].to_numpy()[i] - 1] += 1
print([RegCounts, InfCounts])
total = sum(RegCounts) + sum(InfCounts)
chi2 = 0
for i in range(len(RegCounts)):
    for col in [RegCounts, InfCounts]:
        Eij = (sum(col) * sum([RegCounts[i], InfCounts[i]])) / total
        chi2 += (col[i] - Eij)**2 / Eij
pval = 1 - stats.chi2.cdf(chi2, 4)
print(f"We get a value of {chi2 = }, and a corresponding {pval = }")
```

we obtain the table,

Response	Respondent Type	
	Regular	Influencer
Str Disagree	19	27
Disagree	18	15
Neutral	21	6
Agree	22	18
Str Agree	20	32
Totals	100	98

We have

$$\text{df} = (\text{nrows} - 1) \times (\text{ncols} - 1) = (2 - 1) \times (5 - 1) \\ = 4$$

degrees of freedom here, and we see that >80% of cells in the table have  $\geq 10$  counts, and all cells have  $\geq 5$  counts, so Pearson's  $\chi^2$  is

of cells in the table have  $\geq 10$  counts, and all cells have  $\geq 5$  counts, so Pearson's  $\chi^2$  is a suitable statistic.

From the python code, we get

$$\begin{aligned}\chi^2 &\approx 13.148 \Rightarrow p\text{-val} = P(\chi^2_4 \geq \chi^2) \\ &= P(\chi^2_4 \geq 13.148) \\ &\approx 0.0106\end{aligned}$$

So we obtain significant evidence that there is at least one differing opinion between a population of regular people and influencers.

(i.e. that regular people and influencers somewhat disagree on the removal of 'likes' from this particular social media platform.)

c. Check audio file. To do:

- We assume that the data points were independent of each other, i.e. that no one person "influenced" (haha) another person's answer.
- People with 10000+ followers are representative of a sample of "influencers"
- Pearson's  $\chi^2$  is an accurate test statistic for this data set. We determined that the counts satisfied the minimum requirements to use Pearson's  $\chi^2$ , but ideally all cells would have a count of 10 or more.

d. Firstly, the two-sample t-test relies on the assumption that the data is described by an underlying standard normal distribution.

We don't know if the data can be described by this.

Secondly, our domain of discrete measurements is simply too small to accurately approximate our discrete data with a continuous normal distribution.

#### Q4

We have a two-way ANOVA, with

$$Y_{jki} = \mu + \alpha_j + \beta_k + \delta_{jk} + \epsilon_{jki}$$

with constraints

$$\sum_j \alpha_j = 0; \sum_k \beta_k = 0; \sum_j \delta_{jk} = \sum_k \delta_{jk} = 0$$

a. Under the sum constraints, the joint likelihood function is given by

$$L(\mu, \alpha, \beta, \delta, \sigma^2 | Y) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \exp \left( -\frac{1}{2\sigma^2} \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^{n_{jk}} (Y_{jki} - \mu - \alpha_j - \beta_k - \delta_{jk})^2 \right)$$

We can then maximize the log likelihood:

$$\log L = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^{n_{jk}} (Y_{jki} - \mu - \alpha_j - \beta_k - \delta_{jk})^2 \quad (1)$$

Differentiating (1) w.r.t.  $\mu$  and setting to 0 gives the MLE for  $\mu$ ,

$$\Rightarrow \frac{\partial \log L}{\partial \mu} = 0 = -\frac{1}{\sigma^2} \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^{n_{jk}} (Y_{jki} - \hat{\mu} - \alpha_j - \beta_k - \delta_{jk})$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^{n_{jk}} Y_{jki}$$

$$\Rightarrow \frac{\partial \log L}{\partial \mu} = 0 = -\frac{1}{\sigma^2} \sum_{j=1} \sum_{k=1} \sum_{i=1} (Y_{jki} - \hat{\mu} - \alpha_j - \beta_k - \delta_{jki})$$

$$\Rightarrow N\hat{\mu} = \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} - Kr \sum_{j=1} \alpha_j - Jr \sum_{k=1} \beta_k - \sum_{j=1} \sum_{k=1} \sum_{i=1} \delta_{jki}$$

with the sum constraints, this is then

$$N\hat{\mu} = \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} - 0 - 0 - 0$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} = \bar{Y}_{...}$$

where  $N = J \cdot K \cdot r$

Similarly, differentiating (1) w.r.t.  $\alpha_j$  and setting to 0 gives

$$\frac{\partial \log L}{\partial \alpha_j} = 0 = -\frac{1}{\sigma^2} \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} - \mu - \hat{\alpha}_j - \beta_k - \delta_{jki}$$

$$\Rightarrow Kr\hat{\alpha}_j = \sum_{k=1} \sum_{i=1} Y_{jki} - Kr\mu - r \sum_{k=1} \beta_k - \sum_{k=1} \sum_{i=1} \delta_{jki}$$

Applying the sum constraints gives

$$Kr\hat{\alpha}_j = \sum_{k=1} \sum_{i=1} Y_{jki} - Kr\mu - r - 0 - 0$$

$$\Rightarrow \hat{\alpha}_j = \frac{1}{Kr} \sum_{k=1} \sum_{i=1} Y_{jki} - \mu$$

$$= \bar{Y}_{j...} - \bar{Y}_{...}$$

for  $\beta_k$ , we get

$$\frac{\partial \log L}{\partial \beta_k} = 0 = -\frac{1}{\sigma^2} \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} - \mu - \alpha_j - \hat{\beta}_k - \delta_{jki}$$

$$\Rightarrow Jr\hat{\beta}_k = \sum_{j=1} \sum_{i=1} Y_{jki} - Jr\mu - r \sum_{j=1} \alpha_j - \sum_{j=1} \sum_{i=1} \delta_{jki}$$

Applying sum constraints gives

$$\hat{\beta}_k = \frac{1}{Jr} \sum_{j=1} \sum_{i=1} Y_{jki} - \mu - 0 - 0$$

$$= \bar{Y}_{..k} - \bar{Y}_{...}$$

And finally for  $\delta_{jki}$ ,

$$\frac{\partial \log L}{\partial \delta_{jki}} = 0 = -\frac{1}{\sigma^2} \sum_{j=1} \sum_{k=1} \sum_{i=1} Y_{jki} - \mu - \alpha_j - \beta_k - \hat{\delta}_{jki}$$

$$\Rightarrow JK\hat{\delta}_{jki} = \sum_{i=1} Y_{jki} - JK\mu - JK\alpha_j - JK\beta_k$$

$$\Rightarrow \hat{\delta}_{jki} = \frac{1}{JK} \sum_{i=1} Y_{jki} - \mu - \alpha_j - \beta_k$$

$$= \bar{Y}_{j..k} - \bar{Y}_{...} - \bar{Y}_{j...} - \bar{Y}_{..k}$$

as required.

b. WTS:  $SS_{\text{total}} = SS_A + SS_B + SS_{AB} + SS_{\text{resid}}$

We have that  $SS_{\text{total}} = \sum_{jki} (Y_{jki} - \bar{Y}_{...})^2$ , and

$$Y_{jki} - \bar{Y}_{...} = (Y_{jki} - \bar{Y}_{j..k}) + (\bar{Y}_{j..k} - \bar{Y}_{...}) + (\bar{Y}_{..k} - \bar{Y}_{...}) + (\bar{Y}_{j..k} - \bar{Y}_{j..k} - \bar{Y}_{..k} - \bar{Y}_{...})$$

$$\Rightarrow \sum_{jki} (Y_{jki} - \bar{Y}_{...})^2 = \sum_{jki} [(Y_{jki} - \bar{Y}_{j..k})^2 + (\bar{Y}_{j..k} - \bar{Y}_{...})^2 + (\bar{Y}_{..k} - \bar{Y}_{...})^2 + (\bar{Y}_{j..k} - \bar{Y}_{j..k} - \bar{Y}_{..k} - \bar{Y}_{...})^2]$$

$$+ [Y_{jki} - \bar{Y}_{j..k}] (\bar{Y}_{j..k} - \bar{Y}_{...} + \bar{Y}_{..k} - \bar{Y}_{...} + \bar{Y}_{j..k} - \bar{Y}_{j..k} - \bar{Y}_{..k} - \bar{Y}_{...}) \\ + [\bar{Y}_{j..k} - \bar{Y}_{...}] (Y_{jki} - \bar{Y}_{j..k} + \bar{Y}_{..k} - \bar{Y}_{...} + \bar{Y}_{j..k} - \bar{Y}_{j..k} - \bar{Y}_{..k} - \bar{Y}_{...})$$

$$\begin{aligned}
& + [\bar{Y}_{jki} - \bar{Y}_{jki}] (\bar{Y}_{jk..} - \bar{Y}_{...} + \bar{Y}_{ik..} - \bar{Y}_{...} + \bar{Y}_{ji..} - \bar{Y}_{jk..} - \bar{Y}_{...}) \\
& + [\bar{Y}_{j..} - \bar{Y}_{...}] (\bar{Y}_{jki} - \bar{Y}_{jk..} + \bar{Y}_{ik..} - \bar{Y}_{...} + \bar{Y}_{ji..} - \bar{Y}_{jk..} - \bar{Y}_{ik..} - \bar{Y}_{...}) \\
& + [\bar{Y}_{ik..} - \bar{Y}_{...}] (\bar{Y}_{jki} - \bar{Y}_{jk..} + \bar{Y}_{jk..} - \bar{Y}_{...} + \bar{Y}_{ji..} - \bar{Y}_{jk..} - \bar{Y}_{ik..} - \bar{Y}_{...}) \\
& + [\bar{Y}_{ji..} - \bar{Y}_{jk..} - \bar{Y}_{ik..} - \bar{Y}_{...}] (\bar{Y}_{jki} - \bar{Y}_{jk..} + \bar{Y}_{jk..} - \bar{Y}_{...} + \bar{Y}_{ik..} - \bar{Y}_{...}) \\
= & SS_{\text{resid}} + SS_A + SS_B + SS_{AB} \\
& + \sum_{jki} [(Y_{jki} - \bar{Y}_{jki}) (\bar{Y}_{jk..} - 3\bar{Y}_{...}) + (Y_{j..} - \bar{Y}_{...})(Y_{jki} - 2\bar{Y}_{...}) \\
& + (\bar{Y}_{ik..} - \bar{Y}_{...})(Y_{jki} - 2\bar{Y}_{...}) + (\bar{Y}_{jk..} - \bar{Y}_{jk..} - \bar{Y}_{ik..} - \bar{Y}_{...})(Y_{jki} - \bar{Y}_{jk..} + \bar{Y}_{jk..} + \bar{Y}_{ik..} - 2\bar{Y}_{...})] \\
= & SS_A + SS_B + SS_{AB} + SS_{\text{resid}} \\
& + \sum_{jki} 0 \\
= & SS_A + SS_B + SS_{AB} + SS_{\text{resid}}
\end{aligned}$$

c. We'd expect that the residual sum of squares has a distribution given by

$$\frac{SS_{\text{resid}}}{\sigma^2} \sim \chi^2_{df_{\text{resid}}}$$

as the residuals over the variance should have a standard normal distribution, and so the square of the residuals should have a  $\chi^2$  distribution. The sum of the  $SS_{\text{resid}}/\sigma^2$  should be  $\sim \chi^2_{df_{\text{resid}}}$ , as we would expect the degrees of freedom to be the total df of all data, minus the number of means we are estimating. That is

$$\begin{aligned}
df_{\text{resid}} & = df_{\text{tot}} - \text{no. } \mu_i \\
& = JK(r) - JK \\
& = JK(r-1)
\end{aligned}$$

which, for  $J=K=2$  and  $r=3$ , we get  
 $df_{\text{resid}} = 8$ .