

Custom Tables

This document describes the algorithms used in the Custom Tables procedure.

Pearson's Chi-Square

Notation

The following notation is used for the computation of Pearson's chi-square:

R	Number of rows in the sub-table.
C	Number of columns in the sub-table.
f_{ij}	Sum of case weights in cell (i,j).
r_i	Sum of case weights in i-th row, $r_i = \sum_{j=1}^C f_{ij}$.
c_j	Sum of case weights in j-th column, $c_j = \sum_{i=1}^R f_{ij}$.
W	Sum of all case weights, $W = \sum_{j=1}^C c_j = \sum_{i=1}^R r_i$.
E_{ij}	Expected cell counts.
χ_p^2	Pearson's Chi-Square statistic.
p_{ij}	Population proportion for cell (i,j).
$p_{i.}$	Marginal population proportion for i-th row.
$p_{.j}$	Marginal population proportion for j-th column.
df	Degrees of Freedom.
p	p-value of the chi-square test.

2 Custom Tables

α Significance level supplied by the user.

Conditions and assumptions

- Tests will not be performed on Comperimeter tables or tables with MR variables.
- Chi-square tests are performed on each innermost sub-table of each layer.
- If a scale variable is in the layer, that layer will not be used in analysis.
- The row variable and column variable must be two different categorical variables.
- The contingency table must have at least two non-empty rows and two non-empty columns.
- Non-empty rows and columns do not include subtotals and totals.
- Empty rows and columns are assumed to be structural zeros. Therefore, R and C are the numbers of non-empty rows and columns in the table.
- If weighting is on, cell statistics must include weighted cell counts or weighted simple row/column percents; the analysis will be performed using these weighted cell statistics. If weighting is off, cell statistics must include cell counts or simple row/column percents; the analysis will be unweighted.
- Tests are constructed by using all visible categories. Hiding of categories and showing of user-missing categories are respected.

Statistics

Hypothesis:

$$H_0 : p_{ij} = p_{i.}p_{.j} \quad i = 1, \dots, R \quad \text{and} \quad j = 1, \dots, C \quad \text{v.s.} \quad \text{not } H_0$$

Statistic:

$$\chi_p^2 = \sum_{ij} \frac{(f_{ij} - E_{ij})^2}{E_{ij}}, \text{ where } E_{ij} = \frac{r_i c_j}{W}.$$

Under the null hypothesis, the statistic has a Chi-square distribution with $df=(R-1)(C-1)$ degrees of freedom.

Alternatively, the chi-square statistics and degrees of freedom can be computed as the following,

$$\chi_p^2 = \sum_{E_{ij}>0} \frac{(f_{ij} - E_{ij})^2}{E_{ij}},$$

R = #{ r_i >0 } and C = #{ c_j >0 }.

This avoids scanning for empty rows and columns before computations.

P-value:

$$p = 1 - F(\chi_p^2; df),$$

where $F(x; df)$ is the cumulative distribution function of Chi-square distribution with df degrees of freedom.

The chi-square test is significant if the $p < \alpha$.

Use of case weights:

The case weights (or frequency weights) are supposed to be integers representing number of replications of each case. In chi-square tests, we will only check if the

4 Custom Tables

aggregated cell counts f_{ij} are integers. If not, they will be rounded to nearest integer before computations.

Small sample validity of the test

Pearson's chi-square is a large sample test, it may not be valid when sample size is small. A rule of thumb is to check if there are more than 80% of cells have expected cell counts larger than 5 and expected cell counts are all larger than 1.

Column Proportions Test

Notation

The following notation is used for the computation of Column Proportions Tests:

R	Number of rows in the sub-table.
C	Number of columns in the sub-table.
A_i	i -th category of the row variable.
B_j	j -th category of the column variable.
f_{ij}	Sum of case weights in cell (i,j) .
r_i	Sum of case weights in i -th row, $r_i = \sum_{j=1}^C f_{ij}$.
c_j	Sum of case weights in j -th column, $c_j = \sum_{i=1}^R f_{ij}$.
Z	z -statistic.
χ^2	Chi-Square statistic.
p_{ij}	Column proportion for cell (i,j) , such that $\sum_{i=1}^R p_{ij} = 1$
\hat{p}_{ij}	Estimated column proportion for cell (i,j) .
\hat{p}_{ijk}	Estimate of pooled column proportion of j and k -th column in i -th row.

p	p-value of a test.
p_B	Bonferroni corrected p-value.
α	The significance level supplied by the user.

Conditions and Assumptions

- Tests will not be performed on Comperimeter tables, tables with MR variables, and tables with scale variables in the layer.
- Pairwise tests are performed on each row of all eligible innermost sub-tables within each layer.
- Sub-tables must be categorical variable by categorical variable.
- Number of rows and columns must be larger than or equal to two. i.e. $R \geq 2$ and $C \geq 2$.
- Tests are constructed by using all visible categories excluding totals and sub-totals. Hiding of categories and showing of user-missing categories are respected.
- If weighting is on, cell statistics must include weighted cell counts or weighted simple column percents; a weighted analysis will be performed. If weighting is off, cell statistics requested must include cell counts or simple column percents; an unweighted analysis will be performed.
- A proportion will be discarded if the proportion is equal to zero or one, or the sum of case weights in a category is less than 2, (i.e. $c_j < 2$). If less than two proportions are left after discarding proportions, test will not be performed.

Statistics

Table layout:

	B₁	B₂	...	B_C
A₁	p ₁₁	p ₁₂		p _{1C}
A₂	p ₂₁	p ₂₂		p _{2C}
...
A_R	p _{R1}	p _{R2}	...	p _{RC}

Hypothesis:

Without loss of generality, we will only look at the i-th row of the table. In i-th row, $C*(C-1)/2$ comparisons will be made among $p_{i1}, p_{i2}, \dots, p_{iC}$. The (j,k)th hypothesis will be

$$H_{0jk} : p_{ij} = p_{ik} \text{ v.s. } H_{1jk} : p_{ij} \neq p_{ik}.$$

If some proportions in i-th row do not satisfy condition 7, they will be discarded and C will be replaced by number of remaining proportions. If resulting C is less than or equal to one, no comparison will be made.

Aggregated Statistics:

Column proportions tests are based on the aggregated proportions (\hat{p}_{ij}) and cell counts for each column (c_j). Column proportions are computed using the un-

rounded cell counts $\hat{p}_{ij} = \frac{f_{ij}}{c_j}$ which are equal to the proportions actually displayed

in CTABLE.

Statistics for the (i,j)th comparisons:

$$\text{Pooled proportion: } \hat{p}_{ijk} = \frac{\text{round}(c_j)\hat{p}_{ij} + \text{round}(c_k)\hat{p}_{ik}}{\text{round}(c_j) + \text{round}(c_k)}.$$

$$\text{z statistic: } z = \frac{(\hat{p}_{ij} - \hat{p}_{ik})}{\sqrt{\hat{p}_{ijk}(1 - \hat{p}_{ijk})(\text{round}(c_j)^{-1} + \text{round}(c_k)^{-1})}}.$$

$$\text{p-value: } p = 2[1 - \Phi(|z|)],$$

where $\Phi(z)$ is the CDF of standard normal distribution.

Alternatively, the statistics can be constructed as a chi-square statistic,

$$\chi^2 = z^2,$$

the p-value will now be given by $p = 1 - F(\chi^2; 1)$, where $F(x; df)$ is the CDF of chi-square distribution with df degrees of freedom.

A comparison is significant if $p < \alpha$ (or $p_B < \alpha$, if Bonferroni adjusted).

Bonferroni adjustment:

If Bonferroni adjustment for multiple comparisons is requested, the p-value p will be adjusted by

$$p_B = \min\left(\frac{p * C * (C - 1)}{2}, 1\right)$$

8 Custom Tables

Relationship to Pearson's chi-square tests:

The statistics used in column proportion tests is equivalent to the Pearson's chi-square test on a 2x2 table by taking j and k-th column and collapsing all rows except i-th row. Therefore performing column proportion tests on a 2x2 table will give you the same result as Pearson's chi-square test.

Use of case weights:

The case weights (or frequency weights) are supposed to be integers representing number of replications of each case. In column proportions tests, we will only check if the column marginal C_j 's are integers. If not, they will be rounded to the nearest integer.

Pairwise Comparisons

Notation

The following notation is used for the computation of Column Proportions Tests:

k	Number of categories in the sub-table.
k^*	Number of categories with case weights greater than or equal to 2.
μ_i	Population mean of the i-th category, $i=1,\dots,k$.
x_{ij}	j-th observation in i-th group.
w_{ij}	Case weight of the j-th observation in i-th group.
w_i	Sum of case weights in category i, $i=1,\dots,k$.
\bar{x}_i	Mean of category i, $i=1,\dots,k$.
s_i	Standard deviation of category i, $i=1,\dots,k$.
s_{pp}	Pooled standard deviation of all categories.
W	Total case weights. Sum of rounded w_i 's.
p_B	p-value adjusted by using Bonferroni method.
α	Significance level supplied by the user.

Conditions and Assumptions

- Tests will not be performed Comperimeter tables or tables with MR variables.
- Tests are performed on each innermost sub-tables for each layer.
- Row variable must be a scale variable, possibly nested under or over some categorical variables. Column variable must be categorical.
- If weighting is on, cell statistics must include weighted means; a weighted analysis will be performed using the weighted statistics. If weighting is off, cell statistics must include means, an unweighted analysis will be performed.
- Tests are constructed by using all visible, non-empty categories excluding totals and sub-totals. Hiding of categories and showing of user-missing categories are respected.
- Total case weights in each category must be at least two. Categories not satisfying this assumption are not used. If number of categories satisfying this condition is less than two, no comparisons will be made.
- Variances of all categories are assumed to be equal.
- User and system missing values of scale variables are excluded.

Statistics

All Pairwise Comparisons

Hypotheses:

$$H_{0ij} : \mu_i = \mu_j, \text{ v.s. } H_{lij} : \mu_i \neq \mu_j, \text{ for all } i > j.$$

$$\text{Total number of hypotheses: } \frac{k^*(k^* - 1)}{2}, \text{ (where } k^* = \sum_{i=1}^k I(w_i \geq 2) \text{)}.$$

Aggregated statistics:

The statistics in pairwise comparisons are computed from aggregated category means (\bar{X}_i), sample variances (s_i^2) and sample sizes (w_i), $i=1,\dots,k$. Various quantities used in the comparisons are shown below.

$$\text{Total case weight (sample size): } W = \sum_{i=1}^k \text{round}(w_i)I(w_i \geq 2)$$

$$\text{Mean of i-th category: } \bar{X}_i = \frac{\sum_{j=1}^{n_i} w_{ij} X_{ij}}{w_i}$$

$$\text{Sample variance of i-th category: } s_i^2 = \frac{\sum_{j=1}^{n_i} w_{ij} (X_{ij} - \bar{X}_i)^2}{w_i - 1}$$

$$\text{Pooled variance: } s_{pp}^2 = \frac{\sum_{i=1}^k I(w_i \geq 2)(\text{round}(w_i) - 1)s_i^2}{W - k^*}$$

Statistics for (i,j)th comparisons:

Assuming $w_i \geq 2$ and $w_j \geq 2$,

$$\text{T-statistic, } t_{ij} = \frac{(\bar{X}_i - \bar{X}_j)}{s_{pp} \sqrt{\left(\frac{1}{\text{round}(w_i)} + \frac{1}{\text{round}(w_j)} \right)}}$$

$$\mathbf{P\text{-value } p = 2[1 - F(|t_{ij}|; W - k^*)],}$$

where $F(t; n)$ is the distribution function of t-distribution with n degrees of freedom.

$$\mathbf{Bonferroni\ adjusted\ p\text{-value } p_B = \min\left(\frac{pk^*(k^* - 1)}{2}, 1\right).}$$

A comparison is significant if $p < \alpha$ (or $p_B < \alpha$, if Bonferroni adjustment is used).

Possible computation problems:

From the formula, we can see that comparison can be made as long as S_{pp}^2 is nonzero. If variances of categories with cell count greater than or equal to two are all zero, S_{pp}^2 becomes zero and no comparison can be made.

Use of case weights:

The case weights (or frequency weights) are supposed to be integers representing number of replications of each case. If sum of case weights in each group ($W_i, i=1, \dots, k$) are not integers, they will be rounded to the nearest integers before calculations. Consequently, the total weight W will become the sum of rounded W_i 's.