Model of multi-independent samples





Hypothesis tests for several independent samples

Compare several proportions

Compare mean values of several populations

Compare several proportions

Let X be a binary variable taking two values 0 and 1. Collecting data from that variable under k different conditions we have a sample containing k groups of observations related with the conditions

Let $p_1, p_2, ..., p_k$

be probabilities of appearance of value 1 of variable X under each of the above k conditions.

Hypothesis $H: p_1 = p_2 = ... = p_k$

Alternative Hypothesis

K: there is certain difference between $p_1, p_2, ..., p_k$

Data: Perform a 2xk table of 2 rows and k columns: each column for one group, the 1rst row for value 1, the 2nd row for value 0 of the variable at observations:

Table 1. Observed frequency

| | Group 1 | Group2 | Group k | |
|------------------|------------------------|------------------------|---------------------|-------|
| X = 1 | <i>n</i> ₁₁ | <i>n</i> ₁₂ | n_{1k} | n_1 |
| $\mathbf{X} = 0$ | <i>n</i> ₀₁ | n ₀₂ | n _{0k} | n_0 |
| | n ⁽¹⁾ | n ⁽²⁾ | $n^{(k)}$ | п |

$$n_{1} = n_{11} + n_{12} + \dots + n_{1k} \quad ; \quad n_{0} = n_{01} + n_{02} + \dots + n_{0k}$$
$$n^{(j)} = n_{j1} + n_{j0} \quad ; \quad j = 1, 2, \dots, k \quad ; \quad n = n_{0} + n_{1}$$

Compare several proportions

 If the hypothesis is correct, the proportion of occurrence of 1 estimated commonly to all columns (conditions) is equal to

n_1/n

• The proportion of occurrence of **O** estimated commonly to all columns is equal to

 n_0/n

Perform the table of expected (theoretical) frequencies of the hypothesis:

Group k Group 1 Group2 $n^{(k)}$ n_1 $n^{(1)}*$ $n^{(2)}*$ X = 1 n_1 n n n $n^{(2)}$ n_0 $n^{(k)}$ $n^{(1)}*$ X = 0 n_0 n n $n^{(2)}$ $n^{(1)}$ $n^{(k)}$ n

Table 2. predicted (expected) frequency

Perform the table of the test statistic:





LEMMA. Suppose that hypothesis H is true. Then variable χ^2 has distribution approximate to the Chi-square distribution with (k-1)degrees of freedom $\chi^2_{(k-1)}$.

Density function of Chi – squared distribution



Using Excel to Compute *Chi* – *squared* Distribution

- Excel has two functions for computing cumulative probabilities and x values for <u>any</u> Chi - squared distribution:
 - <u>CHI.DIST</u> is used to compute the cumulative probability given an *x* value, p-value.
 - <u>CHI.INV</u> is used to compute the *x* value given a cumulative probability, critical value.

Method A (p-value):

Step 1. Taking a variable $\chi^2_{(k-1)}$ of Chi-squared distribution with (k-1) degrees of freedom calculate the probability (p-value)

b = $P \{\chi^2_{(k-1)} > \chi^2 \}$.

Step 2. Compare the probability **b** to the given ahead significance level α :

* If $b \ge \alpha \rightarrow$ accept hypothesis H, conclude the all proportions are equal

* If $b < \alpha \rightarrow$ reject hypothesis H, confirm the appearance of some difference between proportions.



Method B. (Critical value)

Looking in Table of Chi-squared distribution to find critical value $\chi^2_{(k-1)}(\alpha)$ of Chi-squared distribution with *k-1* degrees of freedom (α is a given ahead significance level =5%,1% or 0.5%)

Decide

- Reject Hypothesis **H**: = if $\chi^2 \ge \chi^2_{(k-1)}(\alpha)$

- Accept Hypothesis H: = if $\chi^2 < \chi^2_{(k-1)}(\alpha)$

The one-factor analysis of variance (ANOVA) test is an extension of the T-test method that compares the mean values to the multi-independent samples' model: For normally distributed random variables $X^{(1)}, X^{(2)}, ..., X^{(k)}$ with common variance and expectations $\mu_1, \mu_2, ..., \mu_k$,

consider the hypothesis:

$$H: \mu_1 = \mu_2 = ... = \mu_k$$

Analysis of Variance:

$$H_0: \ \mu_1 = \mu_2 = \mu_3 = \cdots = \mu_k$$

 H_a : Not all population means are equal

- If H_0 is rejected, we cannot conclude that *all* population means are different.
- Rejecting H_0 means that at least two population means have different values.

Analysis of Variance:

- Assumptions for Analysis of Variance
 - For each population, the response (dependent) variable is normally distributed.
 - The variance of the response variable, denoted σ^2 , is the same for all of the populations.
 - The observations must be independent.

Consider the sample $\{x_i^{(i)}, j = 1, 2, ..., n_i\}$ of observations in the i-th group, corresponding to the variable $X^{(i)}$, including n_i observation, $i = 1, 2, \dots, k$. $x_i^{(i)} = \mu_i + \varepsilon_{ii} = \mu + \alpha_i + \varepsilon_{ii} ,$ Consider the model where \mathcal{E}_{ij} is the deviation of each observation in the group from the group mean μ_i , and μ the common mean of all observations in the data, α_i is the difference between the group mean μ_i and the overall mean, μ obviously $\alpha_1 + \alpha_2 + \ldots + \alpha_k = 0$

The total volatility in the i-th group, i = 1, 2, ..., k, is equal to $SS_i = \varepsilon_{i1}^2 + \varepsilon_{i2}^2 ... + \varepsilon_{in_i}^2$ and $(\varepsilon_{i1}^2 + \varepsilon_{i2}^2 ... + \varepsilon_{in_i}^2)/(n_i - 1)$ is the adjusted variance of the random variable $X^{(i)}$.

is the adjusted variance of the random variable $X^{(r)}$. $RSS = SS_1 + ... + SS_k = (\varepsilon_{11}^2 + \varepsilon_{12}^2 ... + \varepsilon_{1n_1}^2) + ... + (\varepsilon_{k1}^2 + \varepsilon_{k2}^2 ... + \varepsilon_{kn_k}^2)$

is the sum of all volatility in all groups.

Denote the total number of observations by $n = n_1 + ... + n_k$ We can take RSS / (n - k)

as an estimate of the common variance of all variables

 $X^{(1)}, X^{(2)}, \dots, X^{(k)}$

$$BSS = n_1\alpha_1^2 + n_2\alpha_2^2 + \ldots + n_k\alpha_k^2$$

1

is the total variation between the group means

BSS/(k-1)

reflecting the random variation from group to group of the considered quantity .

Ratio between two variances

$$F = \frac{BSS / (k-1)}{RSS / (n-k)}$$

is the test statistic for hypothesis H comparing mean values.

Lemma. If the hypothesis H is true, then the ratio (test statistic) F has a Fisher-Snedecor distribution with (k-1) and (n-k) degrees of freedom.

With the above lemma, the testing problem has approaches of p-value or critical value.

Density function of Fisher-Snedecor distribution (F)



- a) Compare the value of the test statistic F with the critical value $F_{cr} = F_{\alpha}^{(k-1,n-k)}$, which is the $(1 - \alpha)$ percentile of the Fisher-Snedecor distribution with (k-1) and (n-k) degrees of freedom: -
- If $F \ge F_{cr}$ \rightarrow rejects hypothesis H,
- If $F < F_{cr}$ \rightarrow accepts H.

The F Test



b) With **FS** being a random variable with a Fisher-Snedecor distribution with (k-1) and (n-k) degrees of freedom, calculate the probability of significance

$$p_{sig} = P\{FS > F\}$$

Compare the probability of significance of the test statistic with the significance level α (=5%):

-If $p_{sig} \leq \alpha$ \rightarrow rejects hypothesis H,-

- If $p_{sig} > \alpha$ \Rightarrow accepts H.

The right-tailed probability of Fisher-Snedecor distribution (F)



Using Excel to Compute Fisher - Snedecor Distribution

- Excel has two functions for computing cumulative probabilities and x values for <u>any</u> Fisher - Snedecor distribution:
 - <u>F.DIST</u> is used to compute the cumulative probability given an *x* value, p-value.
 - <u>F.INV</u> is used to compute the *x* value given a cumulative probability, critical value.