

The intra-class correlation coefficient (ICC)

Stat 576
5-1-25

①

Form all possible ordered pairs of items

from the same cluster $(Y_{ij}, Y_{ik}), j \neq k$

How many pairs are there? $\sum_{i \in I} M_i(M_i - 1)$

The ICC is the usual correlation, computed on these ordered pairs

An ICC close to 0 is desirable.

Alternate expression for ICC:

②

Source	SS	df
Between	SSB	N-1
Within	SSW	K-N
Total	SST	K-1

$$ICC = 1 - \frac{K}{K-N} \frac{SSW}{SST} \quad (\text{Note the similarity to } R^2_{adj})$$

Worst case: $SSW = 0 \Rightarrow ICC = 1$

Best case: $SSW = SST \Rightarrow ICC = 1 - \frac{K}{K-N}$
 $= -\frac{N}{K-N} \approx 0$

From last time: $\bar{y}_{clus} = \frac{\hat{t}}{K}$, $\hat{t} = N\bar{E}$ (3)

$$V[\bar{y}_{clus}] = \frac{1}{K^2} N^2 \frac{S_t^2}{n} \left(1 - \frac{n}{N}\right)$$

What happens if K is unknown?

This occurs if you don't know

M_i for the clusters not chosen

Then $\hat{t} = N\bar{E}$ is still good

But $\bar{y}_{clus} = \frac{\hat{t}}{K}$ requires an estimator of K

How can we estimate K ? (4)

$$\hat{K} = N \bar{m}$$

↑ sample mean of the
cluster sizes

$$\text{Then } \bar{y}_{clus} = \frac{\hat{t}}{\hat{K}} = \frac{N\bar{E}}{N\bar{m}} = \frac{\bar{t}}{\bar{m}}$$

This is a ratio estimator, so it is asymptotically unbiased

$$V[\bar{y}_{clus}] \approx \frac{1}{\bar{M}^2} \frac{S_e^2}{n} \left(1 - \frac{n}{N}\right), \quad (5)$$

where $S_e^2 = S_t^2 + B^2 S_m^2 - 2B S_{tm}$

\bar{M} = Average cluster size for population = $\frac{K}{N}$

S_t^2 = population variance of cluster totals

$B = \frac{\bar{T}}{\bar{M}}$, \bar{T} = Average cluster total for population = $\frac{t}{N}$

S_m^2 = population variance of cluster sizes (6)

S_{tm} = population covariance between the cluster totals & cluster sizes

$$\hat{V}[\bar{y}_{clus}] = \frac{1}{\bar{M}^2} \frac{S_e^2}{n} \left(1 - \frac{n}{N}\right),$$

$$S_e^2 = S_t^2 + \bar{y}_{clus}^2 S_m^2 - 2\bar{y}_{clus} S_{tm}$$

2-stage cluster sampling

(7)

1st stage: select n clusters from N
by SRSWOR

2nd stage: select a sample of size m_i
from the M_i items in each cluster,
SRSWOR

$$\text{Bias: } \hat{k} = N\bar{E} = N \cdot \frac{1}{n} \cdot \sum_{i=1}^n t_i$$

But t_i is now unknown

Estimate t_i by $\hat{t}_i = M_i \bar{y}_i$

(8)

$$\text{So } \hat{k} = \frac{N}{n} \sum_{i=1}^n M_i \bar{y}_i$$

Law of Iterated Expectations

$$E(E[X|Y]) = E[X]$$

Application to variances:

$$V[X] = E[X^2] - (E[X])^2$$

$$V[X] = E_1(E_2(X^2)) - [E_1(E_2(X))]^2 \quad (9)$$

↑
Conditional expectation, given what
happened at stage 1

$$= \underbrace{E_1(E_2(X^2)) - E_1((E_2(X))^2)}_{\text{}} + \underbrace{E_1((E_2(X))^2) - [E_1(E_2(X))]^2}_{\text{}}$$

$$= E_1[E_2(X^2) - (E_2(X))^2] + E_1((E_2(X))^2) - [E_1(E_2(X))]^2$$

$$V[X] = E_1(V_2(X)) + V_1(E_2(X)) \quad (10)$$

Next time, we will apply this to \hat{t} .