

A study memo on chi-squared test for categorical data

Contents

1	Chi-squared test for categorical data	2
---	---------------------------------------	---

1 Chi-squared test for categorical data

Proposition 1.1. *Let*

- (S1) (Ω, \mathcal{F}, P) is a probability space.
- (S2) $\{X_i\}_{i=1}^\infty$ is a sequence of N -dimensional vectors of random variables on (Ω, \mathcal{F}, P) .
- (A1) $\{X_i\}_{i=1}^\infty$ distribution converges to $N(0, E_N)$.

then $\{|X_i|^2\}_{i=1}^\infty$ distribution converges to $\chi^2(N)$.

Proof. Let us fix arbitrary $a > 0$.

Let λ be the N -dimensional Lebesgue's measure. By (A1) and $\lambda(\partial B(X, \sqrt{a})) = 0$,

$$\begin{aligned} \mu(\{|X_i|^2 \leq a\}) &= \mu(\{X_i \in \overline{B(X, \sqrt{a})}\}) \\ &\rightarrow N(0, E_N)(\overline{B(X, \sqrt{a})}) \quad (i \rightarrow \infty) \end{aligned} \quad (1)$$

By the definition of chi-squared distribution with degree of free N ,

$$N(0, E_N)(\overline{B(X, a)}) = \chi^2(N)([0, a]) \quad (2)$$

So $\{|X_i|^2\}_{i=1}^\infty$ distribution converges to $\chi^2(N)$. \square

Theorem 1.1. *Let*

- (S1) (Ω, \mathcal{F}, P) is a probability space.
- (S2) $\{X_i\}_{i=1}^\infty$ is a sequence of K -dimensional vectors of random variables on (Ω, \mathcal{F}, P) .
- (S3) $\{\pi_k\}_{k=1}^K \subset (0, 1)$ such that $\sum_{k=1}^K \pi_k = 1$.
- (A1) $P(\{X_{i,k} = 1\}) = 1 \quad (\forall i, \forall k)$.
- (A2) For any k, l such that $k \neq l$, $\{X_{i,k} = 1\} \cup \{X_{i,l} = 1\} = \phi \quad (\forall i)$.
- (S4) $O_{n,k} := \sum_{i=1}^n X_{i,k} \quad (n \in \mathbb{N}, k \in \mathbb{N})$.
- (S5) $E_{n,k} := n\pi_k \quad (n \in \mathbb{N}, k \in \mathbb{N})$.

then

$$Q(n) := \sum_{k=1}^K \frac{(O_{n,k} - E_{n,k})^2}{n\pi_k} \quad (3)$$

distribution converges to $\chi^2(K - 1)$.

Proof. We set

$$Y_{n,k} := \sqrt{n}(\bar{X}_k - \pi_k) \quad (n \in \mathbb{N}, k \in \mathbb{N}) \quad (4)$$

Then

$$Y_{n,K} := -\sum_{k=1}^{K-1} Y_{n,k} \quad (\forall n) \quad (5)$$

and

$$O_{n,k} - E_{n,k} = \sqrt{n}Y_{n,k} \quad (n \in \mathbb{N}, k \in \mathbb{N}) \quad (6)$$

$$Y_n := (Y_{n,1}, \dots, Y_{n,K-1})^T \quad (7)$$

If we set $A := \{a_{i,j}\}_{i,j=1,\dots,K-1}$ by

$$a_{i,j} = \begin{cases} \frac{1}{\pi_i} + \frac{1}{\pi_K} & \text{if } (i = j), \\ \frac{1}{\pi_K} & \text{if } (i \neq j), \end{cases} \quad (8)$$

So

$$Q(n) = Y_n^T A Y_n \quad (n \in \mathbb{N}) \quad (9)$$

and A is a nonnegative definite symmetric matrix.

We set $(K-1)$ -by- $(K-1)$ matrix $\Sigma := \{\sigma_{i,j}\}_{i,j=1,\dots,K-1}$ by $\sigma_{i,j} = \text{cov}(X_{1,i}, X_{1,j})$. Then

$$\sigma_{i,j} = \begin{cases} \pi_i(1 - \pi_i) & \text{if } (i = j), \\ -\pi_i\pi_j & \text{if } (i \neq j), \end{cases} \quad (10)$$

and

$$\sigma_{i,j} = \text{cov}(X_{n,i}, X_{n,j}) \quad (\forall n, \forall i, \forall j) \quad (11)$$

By Proposition 1.2, Σ is positive definite symmetric matrix.

By the central limit theorem (see [3]), $Y_{n,n=1}^\infty$ distribution converges to $N(0, \Sigma)$.

By Proposition 1.1, $\{Q(n)\}_{n=1}^\infty$ distribution converges to $\chi^2(K-1)$. \square

Proposition 1.2. *Let A and B be matrixes in the proof of Theorem 1.1. Then $A^{-1} = \Sigma$*

Proof. For any $i \in \{1, 2, \dots, K-1\}$

$$\begin{aligned} (A\Sigma)_{i,i} &= a_{i,i}\sigma_{i,i} + \sum_{k \neq i} a_{i,k}\sigma_{k,i} \\ &= \left(\frac{1}{\pi_i} + \frac{1}{\pi_K}\right)\pi_i(1 - \pi_i) + \sum_{k \neq i} \frac{1}{\pi_K}(-\pi_i\pi_j) \\ &= (1 - \pi_i) + \pi_i \frac{(1 - \pi_i) - \sum_{k \neq i} \pi_k}{\pi_K} \\ &= 1 \end{aligned} \quad (12)$$

For any $i \in \{1, 2, \dots, K-1\}$ and any $j \in \{1, 2, \dots, K-1\}$ such that $i \neq j$,

$$\begin{aligned} (A\Sigma)_{i,j} &= a_{i,i}\sigma_{i,j} + a_{i,j}\sigma_{j,j} + \sum_{k \neq i,j} a_{i,k}\sigma_{k,i} \\ &= \left(\frac{1}{\pi_i} + \frac{1}{\pi_K}\right)(-\pi_i\pi_j) + \frac{1}{\pi_K}\pi_j(1 - \pi_j) + \sum_{k \neq i,j} \frac{1}{\pi_K}(-\pi_k\pi_j) \\ &= \left(-\pi_j - \frac{\pi_j}{\pi_K}\pi_i\right) + \left(\frac{\pi_j}{\pi_K} - \frac{\pi_j}{\pi_K}\pi_j\right) - \frac{\pi_j}{\pi_K}\sum_{k \neq i,j} \pi_k \\ &= -\pi_j + \frac{\pi_j}{\pi_K} - \frac{\pi_j}{\pi_K}(1 - \pi_K) \\ &= 0 \end{aligned} \quad (13)$$

\square

References

- [1] Tadahisa Funaki, Probability Theory(in Japanese), ISBN-13 978-4254116007.
- [2] Shinichi Kotani, Measure and Probability(in Japanese), ISBN4-00-010634-1.
- [3] A study memo on a proof of the central limit theorem
<https://osmanthus.work/?p=607>
- [4] Tatsuya Kubota, Foundations of modern mathematical statistics(in Japanese), ISBN978-4-320-11166-0.