

# $\chi^2$ 適合度検定の初等的証明

## An Elementary Proof of K.Pearson's Chi-square Test for Goodness of Fit

中嶋真澄

Masumi NAKAJIMA

*Department of Economics*

*International University of Kagoshima*

*Kagoshima 891-0191, JAPAN*

e-mail: nakajima@eco.iuk.ac.jp

### 概要

#### Abstract

We prove here K.Pearson's chi-square test for goodness of fit elementarily.

Key words ; K.Pearson's chi-square test for goodness of fit,  $\chi^2$  distribution.

Mathematics Subject Classification 2010; 62E20, 60E99.

$\chi^2$  適合度検定とは、ドイツの測地学者兼天文学者 Helmer 氏が [3] で発見し、その主要性質を [4] で求めた  $\chi^2$  分布 [2] [12] を使って、K.Pearson(1857-1936, 息子の E.S.Pearson(1895-1980) も統計学者) が、1900 年 [10] で提案した、想定する確率分布とデータが適合しているか否かをある確率で判定する、今日の統計学で広く使用されている検定法の一つで以下のようなものである [1] ( [1] は [6] の後、著者の講義「統計学 II」, 「確率と統計 II」で教科書として使い続けているバランスの取れた良書である。 )。 ( [7] では  $\chi^2$  分布とは呼ばずに、その発見者に敬意を表して Helmer 分布と呼んでいる。 ) :

未知の  $r$  項分布母集団  $Z$  に対して、 $Z$  の値域  $R \subset \mathbf{R}$  を、 $R = \bigcup_{i=1}^r A_i$

( $A_i \cap A_j = \emptyset$  for  $i \neq j$ ) と分解し,  $Z \in A_j$  となる確率  $P\{Z \in A_j\} =: p_j$ , ( $0 \leq p_j \leq 1$ ,  $\sum_{j=1}^r p_j = 1$ ) と想定し (帰無仮説), 実際に測定した  $n$  個の  $Z$  の実現値  $z$  で  $z \in A_j$  となる度数を  $n_j$  とするとき ( $\sum_{j=1}^r n_j = n$ ):

事象	$A_1$	$A_2$	$\cdots$	$A_r$
確率	$p_1$	$p_2$	$\cdots$	$p_r$
頻度	$\frac{n_1}{n}$	$\frac{n_2}{n}$	$\cdots$	$\frac{n_r}{n}$

$$\chi^2 = \sum_{i=1}^r \frac{(n_i - np_i)^2}{np_i}$$

が, この仮説の下で  $n \rightarrow \infty$  のとき, 自由度  $(r-1)$  の  $\chi^2$  分布をすることを利用して  $\chi^2$  検定するものである。

この検定は次の定理を使っている。

**定理 1**  $Z_i$  が  $r$  項分布  $\text{Mul}(n; p_1, p_2, \dots, p_r)$  に従っているとき (このとき  $Z_i \sim B(n, p_i)$ , 即ち各  $Z_i$  は 2 項分布  $B(n, p_i)$  に従っている), もしくは  $(Z_1, Z_2, \dots, Z_r) \sim \text{Mul}(n; p_1, p_2, \dots, p_r)$  のとき ( $r = 2$  のときは,  $\text{Mul}(n; p_1, p_2) = B(n, p_1)$  である),

$$\chi^2 = \sum_{i=1}^r \frac{(Z_i - np_i)^2}{np_i}$$

は,  $n \rightarrow \infty$  のとき, 自由度  $(r-1)$  の  $\chi^2$  分布に従う。(実際の適用では,  $n \rightarrow \infty$  ではなく  $np_i \geq 5$  ( $i = 1, 2, \dots, r$ ) で十分である。)

この定理の言明は簡単であるのにもかかわらず, その証明は難しいせいか, 或いは最近の大学生の学力低下の為か, 最近の殆どの統計学の書物にはその証明が記載されていない。又, 誤った記述がされているものさへある (例えば, 過去私が「統計学 II」の教科書として使った [6] の p.130, 132, 134 では, 「 $n \rightarrow \infty$  では」或いは「近似的に」という語が脱落している)。昨今の新聞紙上にをにぎわした血圧降下剤の臨床試験論文のデータ改竄事件を始め, 統計学を使用している論文等に統計学の使用上の誤用不正が多く見受けられる。このような御時世からも, 判断力のない初心者に対する教科書では, 些細な誤りさへ禁忌である。19 世紀英国の首相 Disraeli が述べ, 今日英国では諺の様になっていると言われる「嘘には 3 種類ある。普通の嘘, 見え透いた嘘, そして統計」が今日の日本で体现されているかのようである。この為か, この定理 1 の証明は大学生の水準の高かつ

た時代の教科書 [5] [9] か、統計学のユーザーである研究者ではなく、統計学自体の研究者となるべき人達が読む限られた専門書にしか現在掲載されていない。このような所からも、上記のような教科書に誤記が現れる理由もあるかのようである。

その肝腎の現在ある証明は、特性関数 (Fourier 変換), Lévy の定理の多次元版多次元正規分布等を使うもので、数学が専門でなければ初等的なものとは云えない。この論文では、このような高度な道具立てを一切使わず、上記定理を初等的に証明する。これは「統計学 II」の講義準備の際考え付いたものである。この程度に初等的なものならば、文系大学での統計学の講義でも使用できると思われる。

定理の証明 次の補題を使う。

#### 補題 1 (de Moivre-Laplace の中心極限定理)

$n$  が十分大きいとき、2 項分布  $B(n, p)$  は、正規分布  $N(np, np(1-p))$  で近似できる。即ち、 $X \sim B(n, p)$  ならば

$$\lim_{n \rightarrow \infty} \frac{X - np}{\sqrt{np(1-p)}} \sim N(0, 1)$$

である。

この補題を証明するには、de Moivre-Stirling の公式を使用するのが普通であるが、そうでない簡単な「証明」もある。[8] [11]

$r$  個の確率変数  $Z_i$  ( $i = 1, 2, \dots, r$ ) を  $r$  項分布:  $\text{Mul}(n; p_1, p_2, \dots, p_r)$  に従っているとすると。即ち、 $Z_i = q_i$  ( $q_i = 0, 1, 2, \dots, n$ ,  $q_1 + \dots + q_r = n$ ,  $i = 1, 2, \dots, r$ ) となる確率を

$$P\{Z_1 = q_1, \dots, Z_r = q_r\} = \frac{n!}{q_1! \dots q_r!} p_1^{q_1} \dots p_r^{q_r}$$

$p_1 + p_2 + \dots + p_r = 1$ ,  $0 \leq p_1, p_2, \dots, p_r \leq 1$  とする。 $(Z_j (j \neq i))$  の値は上記制限内でどのような値でも構わないとして) 周辺分布を考えれば、 $Z_i \sim B(n, p_i)$  である。これは

$$\begin{aligned} P\{Z_i = q_i\} &= \sum_{\substack{q_j, j \neq i \\ \sum_{j=1}^r q_j = n - q_i}} \frac{n!}{q_1! \dots q_i! \dots q_n!} p_1^{q_1} \dots p_i^{q_i} \dots p_n^{q_n} \\ &= \frac{n!}{q_i!} p_i^{q_i} \sum_{\substack{q_j, j \neq i \\ \sum_{j=1}^r q_j = n - q_i}} \frac{1}{q_1! \dots 1 \dots q_n!} p_1^{q_1} \dots 1 \dots p_n^{q_n} \end{aligned}$$

$$\begin{aligned}
&= \frac{n!}{q_i!(n-q_i)!} p_i^{q_i} \sum_{q_j, j \neq i: \sum_{j \neq i} q_j = n - q_i} \frac{(n-q_i)!}{\prod_{j \neq i} q_j!} \prod_{j \neq i} p_j^{q_j} \\
&= \frac{n!}{q_i!(n-q_i)!} p_i^{q_i} (1-p_i)^{n-q_i}
\end{aligned}$$

より明らかである。従って、2項分布の性質から  $Z_i$  の平均、分散はそれぞれ

$$E[Z_i] = np_i, \quad V[Z_i] = np_i(1-p_i)$$

となる。これと同様な計算を  $E[Z_i Z_j]$  ( $i \neq j$ ) に行うと

$$\begin{aligned}
&E[Z_i Z_j] \\
&= \sum_{\substack{q_i, q_j=0 \\ 0 \leq q_i + q_j \leq n}}^n q_i q_j \frac{n!}{q_i! q_j! (n-q_i-q_j)!} p_i^{q_i} p_j^{q_j} \times \\
&\quad \times \sum_{q_k, k \neq i, j: \sum_{k \neq i, j} q_k = n - q_i - q_j} \frac{(n-q_i-q_j)!}{\prod_{k \neq i, j} q_k!} \prod_{k \neq i, j} p_k^{q_k} \\
&= \sum_{\substack{q_i, q_j=0 \\ 0 \leq q_i + q_j \leq n}}^n q_i q_j \frac{n!}{q_i! q_j! (n-q_i-q_j)!} p_i^{q_i} p_j^{q_j} \left( \sum_{k \neq i, j} p_k \right)^{n-q_i-q_j} \\
&= \sum_{\substack{q_i, q_j=0 \\ 0 \leq q_i + q_j \leq n}}^n q_i q_j \frac{n!}{q_i! q_j! (n-q_i-q_j)!} p_i^{q_i} p_j^{q_j} (1-p_i-p_j)^{n-q_i-q_j}
\end{aligned}$$

となる。ここで次の補題を使う。

## 補題 2

$$\begin{aligned}
\sum_{\substack{q_i, q_j=0 \\ 0 \leq q_i + q_j \leq n}}^n q_i q_j \frac{n!}{q_i! q_j! (n-q_i-q_j)!} x^{q_i} y^{q_j} z^{n-q_i-q_j} &= xy \frac{\partial^2}{\partial x \partial y} (x+y+z)^n \\
&= n(n-1)xy(x+y+z)^{n-2}
\end{aligned}$$

証明 明らか。□

この補題より、先の続きは

$$\begin{aligned}
E[Z_i Z_j] &= n(n-1)p_i p_j \{p_i + p_j + (1-p_i-p_j)\}^{n-2} \\
&= n(n-1)p_i p_j
\end{aligned}$$

となる。従って  $Z_1, Z_2, \dots, Z_r$  は独立ではない ( $\sum_{k=1}^r Z_k = n$ )。ここで  $X_i$  を次で定義する。

$$X_i := \frac{Z_i - np_i}{\sqrt{np_i}} \quad (i = 1, 2, \dots, r).$$

$X_i$  は

$$E[X_i] = 0, \quad V[X_i] = E[X_i^2] = E\left[\frac{(Z_i - np_i)^2}{np_i}\right] = \frac{V[Z_i]}{np_i} = 1 - p_i$$

を満たし、束縛条件：

$$\sum_{k=1}^r \sqrt{p_k} X_k = 0$$

を、これ 1 つのみ持つ。これが  $X_1, \dots, X_r$  が独立でない理由である。補題 1 より、

$$\lim_{n \rightarrow \infty} X_i \sim N(0, 1 - p_i) \quad (i = 1, 2, \dots, r)$$

である。又、 $E[Z_i Z_j] = n(n-1)p_i p_j$  を使って  $E[X_i X_j]$  を計算すると

$$\begin{aligned} E[X_i X_j] &= E\left[\frac{Z_i - np_i}{\sqrt{np_i}} \cdot \frac{Z_j - np_j}{\sqrt{np_j}}\right] \\ &= \frac{1}{n\sqrt{p_i p_j}} \{n(n-1)p_i p_j - np_i np_j - np_j np_i + n^2 p_i p_j\} \\ &= -\sqrt{p_i p_j} \end{aligned}$$

である。

次に先の  $X_1, \dots, X_r$  を直交変換  $A = (a_{ij})$  with  $a_{rk} = \sqrt{p_k}$ , ( $k = 1, 2, \dots, r$ ),  ${}^t A A = I$  で  $X'_1, \dots, X'_r$  に次で変換する：

$$X'_j = \sum_{k=1}^r a_{jk} X_k.$$

$j = 1, 2, \dots, r-1$  に対して、 $\sum_{k=1}^r \sqrt{p_k} X_k = 0$  を使って  $X'_j$  を計算すると

$$X'_j = \sum_{k=1}^{r-1} (a_{jk} - a_{jr} \frac{\sqrt{p_k}}{\sqrt{p_r}}) X_k$$

となり、

$$E[X'_j] = \sum_{k=1}^r a_{jk} E[X_k] = \sum_{k=1}^r a_{jk} \cdot 0 = 0 \quad (j = 1, 2, \dots, r),$$

$$\begin{aligned}
V[X'_j] &= \sum_{k=1}^r a_{jk}^2 V[X_k] + \sum_{k \neq l} a_{jk} a_{jl} \text{Cov}[X_k X_l] \\
&= \sum_{k=1}^r a_{jk}^2 V[X_k] + \sum_{k \neq l} a_{jk} a_{jl} E[X_k X_l] \\
&= \sum_{k=1}^r a_{jk}^2 (1 - p_k) - \sum_{k \neq l} a_{jk} a_{jl} \sqrt{p_k p_l} \\
&= \sum_{k=1}^r a_{jk}^2 - \sum_{k=1}^r a_{jk}^2 p_k - \sum_{k \neq l} a_{jk} a_{jl} \sqrt{p_k p_l} \\
&= 1 - \left( \sum_{k=1}^r a_{jk} \sqrt{p_k} \right)^2 \\
&= 1 - \left( \sum_{k=1}^r a_{jk} a_{rk} \right)^2 = 1 - 0 = 1 \text{ for } j = 1, 2, \dots, r-1
\end{aligned}$$

である。ここで  $a_{jk}$  が直交行列の成分であることを使った。  
 $X'_r$  については、

$$\begin{aligned}
V[X'_r] &= \sum_{k=1}^r a_{rk}^2 V[X_k] + \sum_{k \neq l} a_{rk} a_{rl} \text{Cov}[X_k X_l] \\
&= \sum_{k=1}^r a_{rk}^2 (1 - p_k) - \sum_{k \neq l} a_{rk} a_{rl} \sqrt{p_k p_l} \\
&= \sum_{k=1}^r a_{rk}^2 - \sum_{k=1}^r a_{rk}^2 p_k - \sum_{k \neq l} a_{rk} a_{rl} \sqrt{p_k p_l} \\
&= 1 - \left( \sum_{k=1}^r a_{rk} \sqrt{p_k} \right)^2 \\
&= 1 - \left( \sum_{k=1}^r a_{rk} a_{rk} \right)^2 = 1 - 1 = 0,
\end{aligned}$$

従って、 $n \rightarrow \infty$  で  $X'_r, X'_i (i = 1, \dots, r-1)$  は正規分布をしているから確率1で  $\lim_{n \rightarrow \infty} X'_r = 0$  となる。

$$X'_j = \sum_{k=1}^{r-1} (a_{jk} - a_{jr} \frac{\sqrt{p_k}}{\sqrt{p_r}}) X_k$$

により、 $n \rightarrow \infty$  で  $X'_j (j = 1, \dots, r-1)$  は正規分布の再生性から、正規分布に従うが、 $E[X'_j] = 0, V[X'_j] = 1$  for  $j = 1, 2, \dots, r-1$  であったことから、更に

$$\lim_{n \rightarrow \infty} X'_j \sim N(0, 1) \text{ for } j = 1, 2, \dots, r-1$$

となる。

又,  $i \neq j$  に対して

$$\begin{aligned}
 E[X'_i X'_j] &= \sum_{k=1}^r \sum_{l=1}^r a_{ik} a_{jl} E[X_k X_l] \\
 &= - \sum_{k=1}^r \sum_{l=1}^r a_{ik} a_{jl} \sqrt{p_k p_l} \\
 &= - \sum_{k=1}^r \sum_{l=1}^r a_{ik} a_{jl} a_{rk} a_{rl} \\
 &= - \sum_{k=1}^r a_{ik} a_{rk} \sum_{l=1}^r a_{jl} a_{rl} \quad (i, j \neq r) \\
 &= -0 \cdot 0 = 0
 \end{aligned}$$

となるので

補題 3  $X, Y$  が正規分布に従う場合,  $X, Y$  が無相関である事と, 独立である事は同値である。即ち

$$\begin{aligned}
 &In \text{ case of } X \sim N(\mu_X, \sigma_X^2), Y \sim N(\mu_Y, \sigma_Y^2), \\
 &X, Y \text{ are uncorrelated} \Leftrightarrow X, Y \text{ are independent.}
 \end{aligned}$$

より,  $X'_1, \dots, X'_{r-1}$  は独立であり,  $n \rightarrow \infty$  のとき, 以下の補題 4 により

$$X_1'^2 + X_2'^2 + \dots + X_{r-1}'^2 \sim \chi^2(r-1)$$

となる。従って  $A = (a_{ij})$  が直交変換であることと  $X'_r = 0$  a.s. であることから

$$X_1^2 + X_2^2 + \dots + X_r^2 = X_1'^2 + X_2'^2 + \dots + X_{r-1}'^2 \text{ with probability 1}$$

となる。これより

$$X_1^2 + X_2^2 + \dots + X_r^2 = \sum_{i=1}^r \frac{(Z_i - np_i)^2}{np_i}$$

が,  $n \rightarrow \infty$  で自由度  $(r-1)$  の  $\chi^2$  分布に従うこと,  
即ち  $n \rightarrow \infty$  のとき

$$\begin{aligned}
 \sum_{i=1}^r \frac{(Z_i - np_i)^2}{np_i} &= X_1^2 + X_2^2 + \dots + X_r^2 \\
 &= X_1'^2 + X_2'^2 + \dots + X_{r-1}'^2 \sim \chi^2(r-1) \\
 &\text{with probability 1}
 \end{aligned}$$

が証明された。ここで,

補題 4  $X_i \sim N(0, 1)$  ( $i = 1, 2, \dots, n$ ) が独立ならば,

$$X_1^2 + X_2^2 + \dots + X_n^2 \sim \chi^2(n)$$

即ち,  $X_1^2 + X_2^2 + \dots + X_n^2$  は自由度  $n$  の  $\chi^2$  分布に従う。

注 補題 4 を自由度  $n$  の  $\chi^2$  分布の定義 definition としても良い。  
を使った。□

## 参考文献

- [1] 儀我真理子 Giga, M.: 『確率・統計の基礎』ムイスリ出版, 2014, (vii+213)pp.. *Foundation of Probability and Statistics*, Muicuri-Shuppan, Tokyo, 2014.
- [2] 林周二 Hayashi, Sh.: 『統計学講義第2版』丸善, *The Lecture on Statistics*, 2nd ed., Maruzen, Tokyo, 1973, (x+378)pp..
- [3] Helmer, F.R.: *Z. Math. und Phys.*, **214**(1876), 192-218.
- [4] Helmer, F.R.: *Astron. Nachr.*, **88**(1876), 113-132.
- [5] 河田敬義, 丸山文行 Kawada, Y., Maruyama, F.: 『数理統計』裳華房, *Mathematical Statistics*, Shouka-Bo, Tokyo, 1951, (vi+230)pp..
- [6] 松井敬 Matsui, T.: 『統計的推測』共立出版, *Statistical Inference*, Kyouritsu-Shuppan, Tokyo, 2012, (viii+207)pp..
- [7] von Mises, R. (edited and complemented by H. Geiringer.): *Mathematical Theory of Probability and Statistics*, Academic Press, 1964, (xiv+694)pp..
- [8] Neyman, J.: *First Course in Probability and Statistics*, Henry Holt and Company, 1950.  
邦訳: 砂田吉一『ネイマン統計学』白桃書房, Tokyo, 1978, (ix+298)pp..
- [9] 小川潤次郎 Ogawa, J.: 『近代数理統計学序説』恵文堂, *Introduction to Modern Mathematical Statistics*, Keibun-Dou, Osaka, 1954, (iii+396)pp..



- [10] Pearson, K.: On the Criterion that a given system of deviations is such that it can be reasonably supposed to have arisen from random sampling, *Phil. Mag.*, **50**(1900), p.157.
- [11] 田島一郎 Tajima, I.: 二項分布から正規分布へ From Binomical Distribution to Normal Distribution, 数学セミナー, *Seminar for Mathematics*, 1974年3月号, March 1974, p.24-26, 日本評論社, Tokyo.
- [12] 吉田伸生 Yoshida, N. 『確率の基礎から統計へ』 *From Foundation of Probability to Statistics*, 遊星社, Yuusei-sha, Tokyo, 2012, (o+189)pp..

(received 23 August 2020.)