

カイ 2 乗検定：1 次元度数分布

例えば、表 1 の形のデータに対して、目の数の頻度が出現比率がお互いに等しい母集団からのものと考えてよいかどうかというような検定について考える。

表 1 サイコロを 100 回振ったときの目の数の出現頻度

目の数	1	2	3	4	5	6
出現頻度	13	16	23	14	19	15

表 1 を一般的にあらわすと表 2 の形となる。

表 2 n 個の観察データを k 個のカテゴリーに分類して数えたもの。 f_i はカテゴリー i に分類されたものの数。データがカテゴリー i に分類されるものである確率を p_i としたとき、観測度数の期待値は $e_i = np_i$ となる。

	カテゴリー 1	・	・	・	カテゴリー k
観察度数	f_1	・	・	・	f_k
確 率	p_1	・	・	・	p_k

表 2 では、n 個のデータ（観測値）が、k 個のカテゴリーに分類されている。カテゴリー i に分類されたものの観測度数が、 f_i で表わされている。カテゴリー i に分類される確率を p_i とおくと、観測度数の期待値 e_i は次式で与えられる。

$$e_i = np_i$$

カイ 2 乗検定では、この期待値 e_i を、データとして与えられた観測度数 f_i と比べることによって、「確率が p_i で与えられる」という仮説の妥当性を調べる。

すなわち、次式

$$c^2 = \sum_{i=1}^k \frac{(f_i - e_i)^2}{e_i} \quad (1)$$

で与えられる値 c^2 を、自由度 $k - 1$ のカイ 2 乗分布と比較することによって、帰無仮説「カテゴリ i の確率は p_i である」の検定を行う。 f_i と e_i の値の違いが大きいほど、式 (1) で与えられる c^2 の値も大きくなる。

有意水準 a で帰無仮説「カテゴリ i の確率は p_i である」を検定するときは、自由度 $k-1$ のカイ 2 乗分布に従う確率変数 X に対して、次式

$$P(X \geq c_a^2) = a$$

を満たす基準値 c_a^2 を求め、式 (1) により算出される c^2 の値と比べる。式 (1) により与えられる c^2 の値が c_a^2 の値より大きいとき、 e_i と f_i の違いが大き過ぎるとして、有意水準 a で帰無仮説「カテゴリ i の確率は p_i である」は棄却される。

プログラム PChiSqr1S.dpr は、等確率

$$p_1 = p_2 = \dots = p_k = \frac{1}{k}$$

を帰無仮説とするものである。このとき、期待値 e_i は

$$e_i = np_i = \frac{n}{k}$$

となる。

プログラム PChiSqr1S.dpr を実行すると、図 1 のフォームが表示される。

	条件 1
ラベル	ラベル 1
度数	

削除 追加 保存 読出 計算 印刷 終了

図1 起動時のフォーム

StringGrid コンポーネント内に設定するデータは、「追加」あるいは「削除」ボタンのクリックを適当に繰り返すことによりセルの列数をデータ数に合わせる。列を増やすときは、「追加」ボタンをクリックする。「追加」ボタンをクリックすると、アクティブなセルを含む行の下に空白行が挿入される。セルは、クリックされるとアクティブになる。

データの設定は StringGrid 内のすべてのセルに対して行わねばならない。データの設定されていない空白のセルがあってはならない。余分な列は、「削除」ボタンのクリックで削除することができる。削除は、アクティブなセルを含む列が除かれる。

図2はデータを設定した状態である。

	条件 1	条件 2	条件 3	条件 4	条件 5	条件 6
ラベル	1の目	2の目	3の目	4の目	5の目	6の目
度数	13	16	23	14	19	15

削除 追加 保存 読出 計算 印刷 終了

図2 データを設定した状態

「保存」ボタンをクリックすると、設定したデータを保存することができる。「保存」ボタンをクリックすると、まず図3のダイアログボックスが表示される。

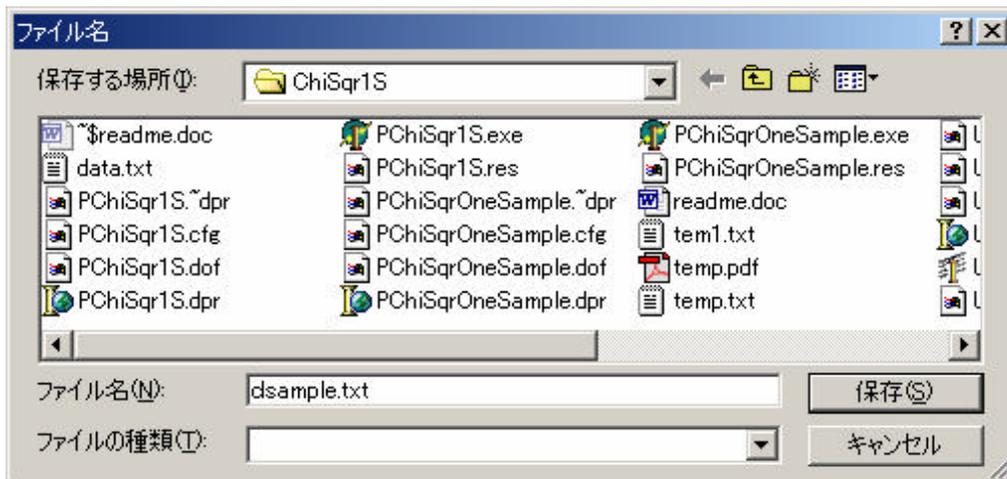


図3 データ保存用ファイル名の設定

ファイルの拡張子は何でもよいが、データはテキストファイルとして保存されるので、図3では拡張子を.txtとしている。データ保存用のファイル名の設定後、図3の「保存」ボタンをクリックすると、設定した名前のファイルにデータが保存される。

保存したデータは、「読出」ボタンのクリックで読み込むことができる。「読出」ボタンをクリックすると、まず読み出すファイルの名前を設定するダイアログボックスが図4のように表示される。

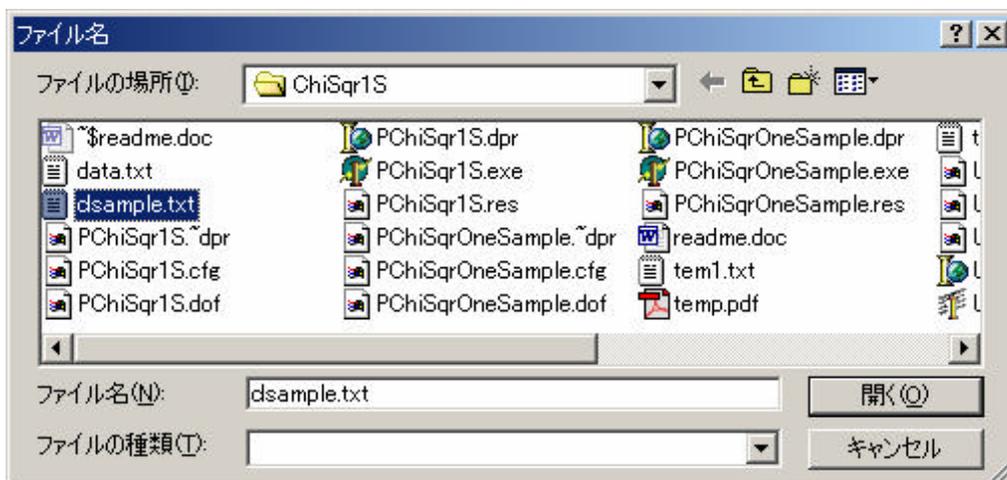


図4 データ読み込みファイル名の設定

ファイル名の設定後、図4の「開く」ボタンをクリックすると、設定したファイルからデータが読み込まれる。

設定されているデータは、「印刷」ボタンのクリックでプリンタに出力することができる。

図2のようにデータを設定した状態で「計算」ボタンをクリックすると、設定されたデ

ータに対するカイ 2 乗の値が算出される。「計算」ボタンのクリックで、まず図 5 のダイアログボックスが表示される。

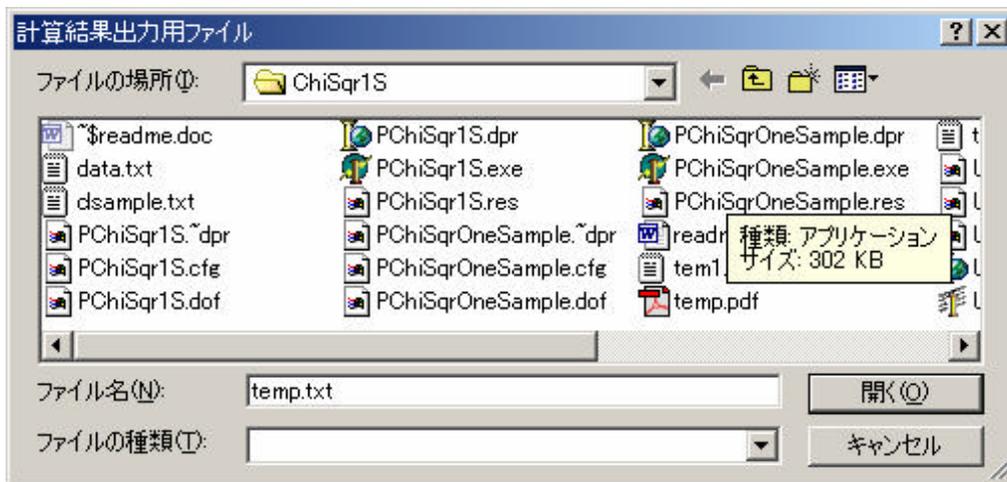


図 5 計算結果出力用ファイル名の設定

図 5 で設定した名前のファイルに、計算結果がテキストファイルとして出力される。名前の設定後、「開く」ボタンをクリックすると計算が始り、計算が終了すると図 6 のフォームになる。

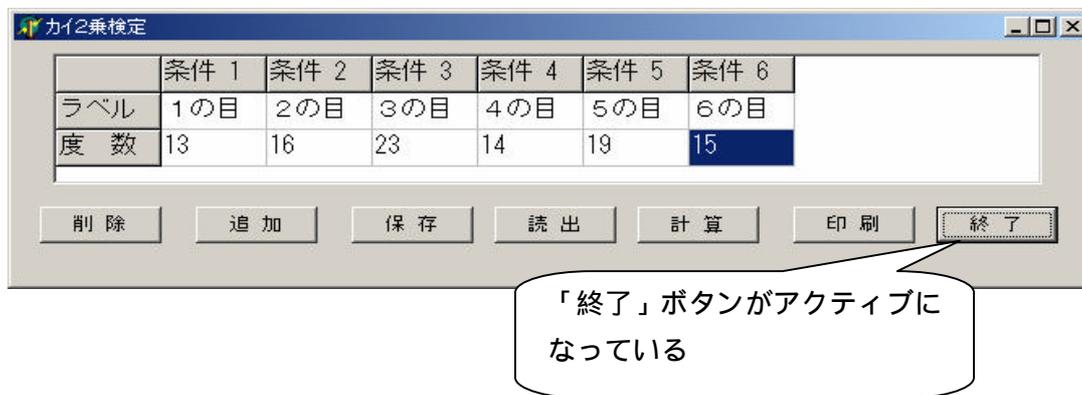


図 6 計算終了時のフォーム

計算が終了すると「終了」ボタンがアクティブになる。「終了」ボタンのクリックでプログラムの実行が終了する。

プログラムの実行終了後、図 5 で設定した名前のファイルを開くとリスト 1 のようにな

っている。

リスト 1 計算結果の出力例

データ	13	16	23	14	19	15
総数 =	100					
e =	16.7					
Chi square =	4.16		df = 5			

リスト 1 に示されているカイ 2 乗の値に対する p 値は

$$p = P(c_5^2 > 4.16) \approx 0.527 > 0.05 = 5\%$$

となる。したがって、リスト 1 の結果の場合、有意水準 5 % で 6 つの目の出現確率には差が認められない。

表 2 の形のデータにカイ 2 乗検定を適用するとき、次の条件が満たされている必要があるとされている (Siegel & Castellan, Jr., 1988)。

- (a) 自由度が 1、すなわちカテゴリー数が 2 の場合は、どちらの期待値 e_i も 5 以上でなければならない。
- (b) 自由度が 2 以上、すなわちカテゴリー数が 3 以上の場合は、期待値 e_i の 20% を越えるものが 5 より小さい、あるいはいずれか 1 つでも 1 より小さいものがあるということがあってはならない。

(a) あるいは (b) の条件が満たされなければならないのは、データ数が十分に多いときに (1) 式で計算される値の分布がカイ 2 乗分布で近似できるからである。(a) の条件が満たされていないときは、2 項分布による検定などを用いる。(b) の条件が満たされていないときは、いくつかのカテゴリーを 1 つにまとめて、 e_i の値が十分に大きくなるよう

にする。カテゴリーの合併は、検定の目的を考えて行う。