

カイ 2 乗検定への理解と数値実験による検証

2012SE119 小嶋佳峰

指導教員：小藤 俊幸

1 はじめに

講義等で学んできたカイ 2 乗検定の原理 [1] は「帰無仮説のもとでは、データから算出されるカイ 2 乗統計量が近似的にカイ 2 乗分布に従う」だが、そもそも、主張自体が分かりにくいと感じたので、本研究では一様乱数を例にして、この原理を数値実験で確かめてみることにした。

2 カイ 2 乗検定

ピアソンのカイ 2 乗検定は「観察された事象の相対的頻度がある頻度分布に従う」という帰無仮説を検定するものである。カイ 2 乗検定のうち最も広く用いられるピアソンのカイ 2 乗検定は適合度検定と独立性検定に用いられる。ここでは適合度の検定について例を使ってまとめる。

K 農事試験場における'えんどう豆'の交配実験結果は次のようだ

種類	円形黄色	円形緑色	角形黄色	角形緑色	計
実測個数	271	78	95	20	464

メンデルの遺伝の法則によれば、これらの個数の比率は、9:3:3:1 となる。上の実験結果は、メンデルの法則に適合しているといえるのか？仮にこの法則どおりに比率がキツカリ 9:3:3:1 になっているのならばそれぞれの個数は総数 464 個の

$$\frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16}$$

になるわけなので、

$$261, 87, 87, 29(\text{個})$$

という数になるはずである。これらの数字

実測個数：	271	78	95	20
理論個数：	261	87	87	29

を比べて、一致してないので「メンデルの法則に適合していない」となるわけではない。この実測個数はあくまでも標本値である。実測値と理論値は一致してはいないがよく似た値である。問題はこれがどの程度近いかということだ。そこでえんどう豆の個数について、実測値と理論値の差をとると、

$$271 - 261 \quad 78 - 87 \quad 95 - 87 \quad 20 - 29$$

これらの差の合計が小さいほどズレも小さいわけだが、このまま合計するとプラス・マイナスが打ち消しあって

$$(271 - 261) + (78 - 87) + (95 - 87) + (20 - 29) = 0$$

のように和は 0 になってしまう。そこでそれぞれの「2 乗」の合計

$$(271 - 261)^2 + (78 - 87)^2 + (95 - 87)^2 + (20 - 29)^2$$

を考えると 0 となる上に、2 乗が実測値と理論値の差を強調できる。しかしこのままではまだ不十分なので、豆一粒あたりのズレを作り、その合計を考える。

$$\frac{(271 - 261)^2}{261} + \frac{(78 - 87)^2}{87} + \frac{(95 - 87)^2}{87} + \frac{(20 - 29)^2}{29} = 4.84$$

261, 87, 87, 29 の時、近似的に自由度 k-1 の χ^2 分布に従うので、有意水準を 0.05 とすると

$$\chi_{k-1}^2(0.05) = \chi_3^2(0.05) = 7.81$$

従って今回のえんどう豆の豆一粒あたりの実測値と理論値の「近さ」は 4.84 で棄却域 7.81 に属さないので H_0 は棄却されない。交配結果はメンデルの法則に適合していないとはいえない結果になった。この種の検定を適合度の検定という。

3 カイ 2 乗検定の数値実験

カイ 2 乗検定の適合度の検定を基に数値実験を行う。メンデルの法則に習った乱数を発生し、そこで発生した値がカイ 2 乗分布に従うかどうかを実験していく。メンデルの遺伝の法則によれば、個数の比率は、9 : 3 : 3 : 1 となるので 0 から 1 の乱数を発生させてそれを if 文を用いて分類分けをすれば 9 : 3 : 3 : 1 の割合で発生する無規則の値が出てくる。

これを M 回繰り返すことで、比率 9 : 3 : 3 : 1 となる M 個の疑似えんどう豆が発生させられる。この M 個のえんどう豆にカイ 2 乗検定を行う。そこで行うカイ 2 乗検定を L 回繰り返す、データを集計する。

最後に配列を用いてカウントし、数値実験を行えるように値を整頓する。これが今回の数値実験のプログラム全体である。

```
#include <stdio.h>
#include <math.h>
#include <stdlib.h>
#define M 1000
#define L 10000

main()
{
    int i, j, p1=0, p2=0, p3=0, p4=0;
    double a, b, kai, l=0, k=0, delta;

    int j1, l1, kai2[101], n;

    srand48(76); //乱数初期化

    for(j1=0; j1<=100; j1++) //配列を初期化
    {
        kai2[j1]=0;
    }

    for(j=0; j<L; j++){ //j を L 回繰り返す
        p1=0, p2=0, p3=0, p4=0;
        for(i=0; i<M; i++){ //i を M 回繰り返す
            b = drand48(); //乱数発生
            if (0.0 <= b && b <= 9.0 / 16.0) { //9:3:3:1
                なるよう場合分け
                p1++;
            }
            else if (9.0 / 16.0 < b && b <= 3.0 / 4.0) {
```

```

        p2++;
    }
    else if (3.0/4.0 < b&& b <= 15.0 /16.0){
        p3++;
    }
    else
    {
        p4++;
    }
}

kai =
((p1-(M*9.0/16.0))*(p1-(M*9.0/16.0)))/(M*9.0/16.0)
+((p2-(M*3.0/16.0))*(p2-(M*3.0/16.0)))/(M*3.0/16.0)
+((p3-(M*3.0/16.0))*(p3-(M*3.0/16.0)))/(M*3.0/16.0)
+((p4-(M*1.0/16.0))*(p4-(M*1.0/16.0)))/(M*1.0/16.0);

//      printf("カイ 2 乗検定:%lf\n",kai);
n=floor(10.0*kai);
//      printf("n:%d\n",n);
if(n<=100)
{
    kai2[n]=kai2[n]+1;
}

}

for(j1=0;j1<=100;j1++){

    delta = (10.0 * kai2[j1])/L ;
    printf("%f\t%f\n",0.1*j1+0.05,delta);

}
return(0);
}

```

また数値実験との比較のための自由度3のカイ2乗分布のプログラムも作成した。

```

#include <stdio.h>
#include <math.h>
#define Pi 3.141592653589793

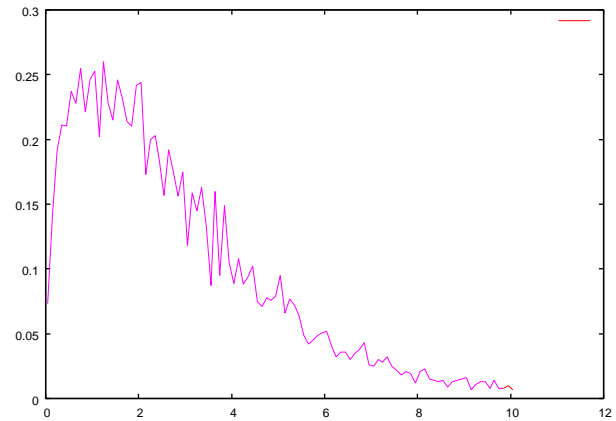
main()
{
    int j1;
    double x,y;

    for(j1=0;j1<=100;j1++)
    {
        x = 0.1 * j1 + 0.05;
        y = (1.0/sqrt(2*Pi))*pow(x,1.0/2.0)*exp(-x/2.0);
        printf("%f\t%f\n",x,y);
    }

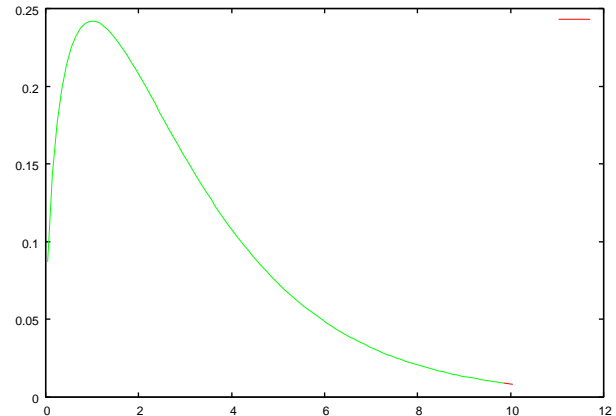
    return(0);
}

```

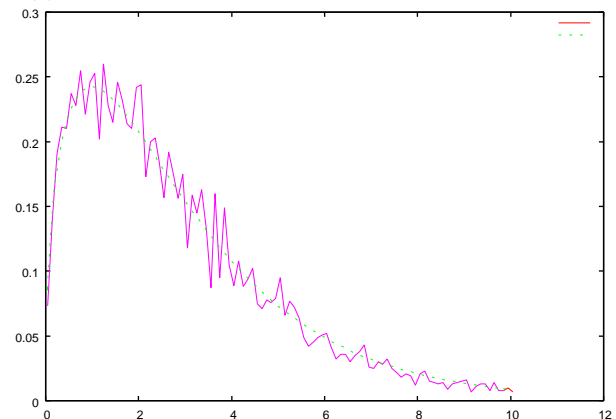
これらを実行し、出てきた値のグラフを比較していく。



今回の数値実験の結果の出力



自由度3のカイ2乗分布



二つのグラフを重ねたグラフ

ここまでを本研究の数値実験とする。

4 おわりに

結果的に数値実験の値がカイ2乗分布に近い値となり、「帰無仮説のもとでは、データから算出されるカイ2乗統計量が近似的にカイ2乗分布に従う」が実証できた。乱数の発生回数を増やせば増やすほどカイ2分布に近づく形のデータがとれたので間違いの無いプログラムの作成ができたと思う。今回はメンデルの法則に沿ったプログラムの作成であったが、値を変更するだけで簡単に他の数値実験も行えるであろう。また応用を加えることで世の中のほとんどの確率に対してのシミュレートが行えるので学生時代の集大成として本研究に取り組みとても良かった。

参考文献

- [1] Ya.G. シナイ (森真訳) 『シナイ確率論入門コース』、1993年
- [2] 小寺平治：『ゼロから学ぶ統計解析』、講談社、2002年