

会資料(2001.2.24)

ベータ 2 項分布の当てはめ

統計数理研

研究所 鈴木 義一郎

§ 1 出生児性比と 2 項分布・ベータ 2 項分布の当てはめ

男児の出生数が女児の出生数に比べて、わずかに多いことはよく知られた統計的事実である。Fisher [1]では、ちょうど 8 人の子供をもつ家庭 53680 組について調べた Geissler のデータについて、男児の生れる確率を p とし、2 項分布 $B(8, p)$ に従うものかどうかを議論している。

まず、表 1 から直接、平均と分散の値を求めてみると

$$\begin{aligned} m &= 221023/53680 = 4.117418 \\ v &= 1021023/53680 - 4.117418^2 \\ &= 2.074167 \end{aligned}$$

のように算出される。一方、2 項分布を想定したときには

$$m = np = 8 \cdot p, \quad v = npq = 8 \cdot p(1-p)$$

であるから、 p についての推定値

$$P = 4.117418/8 = 0.514677$$

を用いて v (分散) の値を求めてみると

$$v = 8 \cdot P(1-P) = 8 \times 0.514677 \times 0.485323 = 1.998277$$

のように算出される。データから直接計算した分散の値より

$$2.074167 - 1.998277 = 0.075890 (3.798\%)$$

だけ大きい。

ところで、この分散に対する標準誤差の平方は、次のように算出される。

$$\begin{aligned} & \frac{2npq}{N-1} + \frac{npq(1-6pq)}{N} = \frac{2 \cdot 8 \cdot 0.514677 \cdot 0.485323}{53679} \\ & \quad + \frac{8 \cdot 0.514677 \cdot 0.485323(1-6 \cdot 0.514677 \cdot 0.485323)}{53680} \end{aligned}$$

5323)

$$\begin{aligned} & + \frac{2 \cdot 1.998277}{53679} + \frac{1.998277 \cdot (0.498708)}{53680} \\ & = \frac{53679}{53680} + \frac{53680}{53680} \\ & = 0.000055888 = (0.007476)^2 \end{aligned}$$

つまり、データから直接計算した分散の値は、2 項分布を仮定した

ときの分散より、その標準誤差の10倍も大きいことが分かる。実際、表2の結果をみても、男だけとか女だけといった家庭の数値が2項分布を仮定した理論度数をかなり上まわっている。さらに、男児が5または3という家庭の数が、逆に少なすぎる傾向もみられる。変動が大きすぎる原因の1つとして、多生児の影響が考えられる。いま、8人の子供の中に、“1卵性”双生児が丁度1組だけいるという家庭(A)での男児の数をXとすると、Xは $B(6, p)$ に従う変数Yと、 $B(1, p)$ に従う変数Zの2倍の和で表すことができる。したがってXの分散は

$$V\{X\} = V\{Y\} + V\{2Z\} = 6pq + 22 \times pq = 10pq$$

となる。

そこで、Aタイプの家庭の比率を $\frac{1}{4}$ として、それ以外の家庭(1 - $\frac{1}{4}$)での男児の数は2項分布 $B(8, p)$ に従うと考えてみよう。ここで2卵性双生児については、同性と異性になる可能性はほぼ等しいものと考えられるし、2組以上の1卵性双生児がいたり、双生児以外の多生児がいる家庭については、無視して考える。すると、全体としての男児の数の分散は

$$\frac{1}{4} \cdot 10pq + (1 - \frac{1}{4}) \cdot 8pq = (1 + \frac{1}{4})8pq$$

のように表すことができる。ところで、このデータの例のように3.8%もの分散の増加があるということから

$$\frac{1}{4} = 0.0382$$

つまり、 $\frac{1}{4} = 0.152$ もの比率でAタイプの家庭が存在しなければならぬことになる。

例えば、1000回の出生のうち12回は双生児であり、さらに双生児のうち1/4は“1卵性”であるといったデータがある。従って1000回の出生のうち3回が1卵性であるから、結局、125(=1000/8)組のうち3組(=2.4%)の家庭が“1卵性双生児”をもつという勘定になる。結局、このデータでの分散の大きさは、1卵性双生児の影響だけで説明するには無理があるということになる。

ところで、分散が2項分布のものより大きい場合には、ベータ2項分布を当てはめるほうがよさそうである。この分布の確率密度関数や平均、分散などが

表2 2項分布(g)とベータ2項分布(h)の適合度

合度

	男児数 期待度数	観測度数	期待度数	偏差		
	(x)	(f)	(g)	(d=f-g)	d^2/g	$(f-h)^2/h$
h	0	215	165.22	+49.78	14.998	193.03
						2.501

1	1485	1401.69	+83.31	4.952	1518.18	0.725
2	5331	5202.65	+128.35	3.166	5334.89	0.003
3	10649	11034.65	-385.65	13.478	10940.00	7.741
4	14959	14627.60	+331.40	7.508	14319.30	28.581
5	11929	12409.87	-480.87	18.633	12250.10	8.416
6	6678	6580.24	+97.76	1.452	6689.27	0.019
7	2092	1993.78	+98.22	4.839	2131.72	0.740
8	342	264.30	+77.70	22.843	303.54	4.873
SUM	53680	53680.00		91.869	53680.00	53.599

$(k;n,a,b) = nCk B(k+a, n-k+b) / B(a,b)$, $(k=0,1,2,\dots, \dots, n)$

$$M = na / (a+b) , V = nab(a+b+n) / (a+b)^2(a+b+1)$$

のように与えられる。与えられた平均と分散の値から、 a , b を推定すると

$$a = 94.35 , b = 88.97$$

といった値が算出される。表 2 には、 $h(k) = (k; 8, 94.35, 88.97)$ を当てはめたときの期待度数も示してある。当てはまりがかなり改良されていることが分かる。

観測値から算出される平均と分散を与えて、ベータ 2 項分布の a , b を推定する関数を次のように定義する：

```
est_ab=:4 :0
```

```
NB. estimate a,b of (k;n,a,b) input n(left);M,V(right)
```

```
a=. (m*m*q=.1 p) -(p=. (m=. {.y.}%x.) *v=. {:y.
```

```
a, (q%p) *a=. a%v m*q
```

```
)
```

```
8 est_ab 4.117418 2.074167
```

```
94.3494 88.9682
```

ベータ 2 項分布 $(k; n, a, b)$ の確率関数は

$(k;n,a,b) = nCk B(k+a, n-k+b) / B(a,b)$, $(k=0,1,2,\dots, \dots, n)$

$$= nCk \cdot \frac{(a+k)!}{(a)!} \cdot \frac{(b+n-k)!}{(b)!} \cdot \frac{(a+b)!}{(a+b+n)!}$$

のように表現できる。そこで、 $(a+k)! / (a)!$ というガンマ関数の比を(ガンマ関数の原始関数「!」を使わないで)算出する関数を次のように定義する。

```
ratio=:4 :0
```

```
NB. calculate (x.+y.) / (y.) not using gamma function(!@< :)
```

```
if. x.=0 do. r=.1 erase. r=.a.<:x.+y.
```

```
while. a>y. do. r=.x:r*a=.a-1 end.
```

```
end.
```

```
)
```

```
bbfn=:4 :0
```

```
NB. beta binomial probability function
```

```
r=. ((|.k) ratio"0{:y.)* (k=.i.x.+1) ratio"0{.y.
```

```
(k!x.) * r%x.ratio+/y.
```

```
)
```

```
0.4": 8 bbfn 94.35 88.97
```

```
0.0036 0.0283 0.0994 0.2038 0.2668 0.2282 0.1246 0.0397 0.0057
```

```
0.2": 53580*8 bbfn 94.35 88.97
```

192.67 1515.35 5324.95 10919.63 14292.60 12227.26 6676.81 212
7.75 302.98

§ 2 ベータ 2 項分布の変形

大相撲の幕内力士は、毎場所 15 人の相手と対戦する。表 3 の $N(x)$ の行に示したのは、平成 6 年 1 月から平成 7 年 7 月までの 10 場所について、途中休場などしなかった力士の勝星数の分布である。どんな相手にも勝率 p で戦える力士の勝星数の確率は、2 項分布 $B(15, p)$ で与えられる。 $A(x)$ の行には、平均勝率 P の値が 0.5 という場合の数値を示してみた(観測データと合わせるために、2 項分布の確率に 386 という数値を掛けて丸めたもの)。

実際のデータと比べてみて、かなりずれている。そもそも、幕内の全ての力士の勝率 P が同じ 0.5 であるとは考えられない。優勝を争うような力士の p は 0.7 以上かもしれないし、反対に不調な下位力士の p は 0.3 程度かもしれない。

そこで、 p の値を“混ぜこぜ”にした「複合 2 項分布」というモデルを考えてやればよい。特に、2 項分布のパラメータ p をベータ分布 $(t; a, b)$ で加重平均したものがベータ 2 項分布 $B(k; 15, a, b)$ である。与えられたデータの平均と分散の値から、 a, b を推定すると

$$a = 9.59, \quad b = 9.47$$

と算出される。そこで、

$$h(k) = h(k; 15, 9.59, 9.47)$$

というモデルを当てはめた結果が、表 3 の $B(x)$ の行に示してある。このモデルにしたら、大分、実際のデータに近くなったが、まだ 7 勝、8 勝の辺りではかなりの“ずれ”が見られる。

ところで 8 勝と 7 勝とでは、次の場所の番付の上下に雲泥の差があり、勝越しをかけた一番は是が非でも勝ちたいと奮起せざるを得ない。したがって、八百長の有無はともかくも、8 勝力士の多くなるのは当然すぎる結果かもしれない。

表 4 に、14 日で 7 勝 7 敗だった力士の数と、千秋楽に勝って 8 勝を達成した力士の数を示してみた。10 場所全体で見ると、7 割 5 分以上の高率を示している。

そこで、14日目までの勝星数(Y)にはベータ2項分布 $g(y) = B(y; 14, a, b)$ を当てはめる。そして、千秋楽での勝星数(Z)の分布としては、14日目でYの値が7である力士の勝率には、“奮起度”がプラスされるものと想定して、

$$\Pr\{Z=1\} = p$$

$$\Pr\{Z=1 | Y=7\} = p +$$

$$\Pr\{Z=1 | Y \neq 7\} = p - \frac{g(7)}{1 - g(7)}$$

といったように与える。この $\frac{g(7)}{1 - g(7)}$ というパラメータは、ZのYに対する従属性の度合いを示すもので、ある種の“dependent outlier”と解釈できる。

次に、15日間での勝星数($X = Y + Z$)の分布を $h(x)$ とすると

$$h(7) = \mu(7) + \frac{d(7)}{2}, \quad h(8) = \mu(8) + \frac{d(8)}{2}$$

$$h(x) = \mu(x) + \frac{d(x)}{2} \quad (x = 7, 8)$$

のように表現できる。ここで

$\mu(x) = \{g(x) + g(x-1)\}/2$, $d(x) = \{g(x) - g(x-1)\}$, [$g(-1) = g(15) = 0$]である。さらに $\frac{g(7)}{1 - g(7)}$ というパラメータも

$$= \frac{h(8) - h(7)}{2[1 - g(7)] + g(6) + g(8)}$$

$$= \frac{h(8) - h(7)}{g(7)\{2 + [g(6) + g(8)]/[1 - g(7)]\}}$$

のように表される。Yの期待値を7と想定すると

$$E\{YZ\} = \sum_y y \Pr\{Z=1, Y=y\} = \sum_y yg(y) \Pr\{Z=1 | Y=y\}$$

$$= (p - \frac{g(7)}{1 - g(7)}) \sum_y yg(y) + (p + \frac{g(7)}{1 - g(7)}) 7g(7) = p \sum_y yg(y) + \frac{7g(7)}{1 - g(7)}$$

$$= E\{Z\}E\{Y\} + \frac{7g(7)}{1 - g(7)} = E\{Z\}E\{Y\}$$

となることから、YとZは無相関になる。したがって

$$V\{Y\} = V\{X\} - V\{Z\} = V\{X\} - p(1-p)$$

という関係より、 $p = 1/2$ と想定してYの分散も推定できるから、 $g(y)$ の分布が特定できることになり、結局、 $\frac{g(7)}{1 - g(7)}$ の値も推定できる([2]参照)。

このように、“dependent outlier”を加えた変形モデルを用いて導出した結果が表3に示したC(x)の行の数値である。B(x)と比べると、7勝、8勝の辺で改良されている。この数値を、実際のデータに重ね合わせてみたのが図1で、ほどほどには適合していることが分かる。

上述のような変形ベータ 2 項分布のモデルを当てはめる関数を次のように定義する。ここで出力結果の最初の値は、従属性のパラメータである の推定値で、以下が当てはめた確率関数の値である。

```
sumou=:3 :0
NB.modified beta binomial model for OZUMOU
v=./f**x +/(x=.i.16)*f=.(%/+)y.
b=.m%1 -m=.7{g=.14 bbfm 14 est_ab 7,v -0.25
c=.b*a=.( -/8 7{f)%m*2+(+/6 8{g)%1 -m
h=.( -:g1+g2)+c*d=(g2=.g,0) -g1=.0,g
a,h+t,(( -c),c),t=.7$0
)
F1=.sumou f1=.0 0 6 13 24 38 53 37 102 54 10 17 17 7 5 3
({.F1);0.1":386*}.F1
```

```
0.4715775 0 2 6 15 40 51 25 90 49 37 24 13 5 2 0
```

```
+/ 0 2 6 15 40 51 25 90 49 37 24 13 5 2 0
```

385

【整数に丸めてしまったので、合計が 386 にならない】

```
F2=.sumou f2=.0 2 8 25 37 80 71 48 150 93 46 23 13 19 9 6
({.F2);0.1":626*}.F2
```

```
0.474372 1 4 12 25 44 64 81 39 141 78 61 41 23 10 3 0
```

また図 2 には、昭和 51 年から 53 年までの 18 場所についての勝星数の分布と、それに当てはめたモデルの結果も示してみた([2])。この時代は、輪島・北の湖そして三重海・二代目若の花が横綱で、それに先代の貴ノ花や旭国といった個性派大関が活躍していた時代である。やはり 7 勝の“落ちこみ”は、現在のパターンと同じような傾向を見せている。しかし、10 勝近辺での“落ちこみ”は今ほど顕著でない。つまり、優勝にからむ力士が大勢いる限り、10 勝力士が極端に少なくなることはなく、したがってここで示したモデルでほぼ説明がつくということにもなる。

【参考文献】

[1] Fisher, R.A. (1970) : 研究者のための統計的方法(遠藤・鍋谷 訳 森北出版)

(Statistical Methods for Research Workers)

[2] 鈴木義一郎(1979) : ある種の Dependent Outlier を含む確率 模型

統計数理研究所彙報 26-1, 23 -

32.