

医療情報によって患者個人が特定される リスクの確率評価

野田 明男

(総合人間科学講座・数学)

Probability Estimation of the Risk of Specifying an Individual Patient by Means of Medical Information

Akio NODA

(Integrated Human Sciences · Mathematics)

Abstract: Let us consider a group of N patients who have been treated in some hospital (N is taken to be 150000 in the case of the hospital affiliated to our university). By using medical information available there, one can classify all patients into a large number of cells. If there exists some cell including only a small number of patients, then each one in that cell has a possibility of being specified as an individual, which might disturb his privacy. The author knows such a risk and its importance, thanks to Professor Kimura in the department of Medical Informatics of our university.

The purpose of the present paper is to establish a probability estimation of the risk mentioned above, in a simple setting of generalized birthday problem ([1]). Noting that the joint probability distribution of the frequencies x_j ($j = 1, 2, \dots, R$) placed in the j -th cell is given by the multinomial distribution, our task is to compute $\alpha = P(x_j < 5 \text{ for some } j)$ which should be very small (less than 1% here). Appealing to the

Poisson approximation ([3]), we arrive at an excellent estimate $\alpha \sim \sum_{j=1}^R \gamma(\lambda_j)$, where $\lambda_j = E(x_j)$ and

$\gamma(\lambda) = \sum_{k=0}^4 \frac{\lambda^k}{k!} e^{-\lambda}$. By virtue of this result, we are led to construct a variety of examples $\{\lambda_j\}$ such that

$\left| \sum_{j=1}^R \gamma(\lambda_j) - 1\% \right| < 0.001\%$, which would help us to investigate complicated real data sets of medical information, collected in our university and also in other universities.

Key words: Generalized birthday problem, Multinomial distribution, Poisson approximation, Probability estimation of risk.

§ 1. 問題の陳述とポアソン近似

ある病院の患者集団 N 人(本学付属病院の場合はおおよそ15万人)を対象にして、種々の医療情報を組み合わせれば、分類をかなり細かくすることができる。そしてある分類項目に属する患者数 x が、限界値 c (ここでは $c=5$ を採用する) より小さくなれば個人が特定されてしまう事態になりかねない。患者のプライバシー保護の立場から、医療情報を含む医学研究の結果公表には、なるべく緩やかな分類を考案するなどして、慎重に対処することが求められている。

さて、対象とする N 人の集団を Ω と書き、1番目の分類基準(例えば糖尿病の有無)を A_i ($1 \leq i \leq r_1$)、2番目の基準(例えば血圧の程度)を B_i ($1 \leq i \leq r_2$)、 \dots 、 h 番目の基準(例えばある薬剤の用量)を H_i ($1 \leq i \leq r_h$) で表す。このとき、細分される分類項目の最小単位(以下セルという)は $A_i \cap B_i \cap \dots \cap H_i$ であり、 Ω は全部で $r_1 r_2 \dots r_h (=R$ とおく) 通りあるセルに、 h 元分割表の形で分割される。

ところで、 h 種類の分類基準間の統計的な関連に着目し、データ分析を実行することが重要であるのは言うまでもない。しかしながら、複雑多岐にわたる医療情報を実際に精査し、研究を進めるのは、医学に不案内な著者には荷が重く、この小論では断念せざるを得ぬ。

こうして、多元分割表の構造そのものは捨てて、 R 通りのセルを単純に一列に並べる(一般化された誕生日問題[1]参照)。そして、各セルに入る確率に応じて、 N 人を割り振る状況を想定する。すなわち、ランダムに選ばれた患者が j 番目のセルに入る確率を p_j ($p_j \geq 0, \sum_{j=1}^R p_j = 1$) とし、 N 人の患者すべてを分類したとき、 j 番目のセルに入った人数が x_j ($x_j \geq 0, \sum_{j=1}^R x_j = N$) となる確率は、多項分布と呼ばれる次式で与えられる。

$$(1-1) \quad P(x_1, \dots, x_R, p_1, \dots, p_R) = \frac{N!}{x_1! \dots x_R!} p_1^{x_1} \dots p_R^{x_R}$$

このとき、個人が特定される危険率は、 $x_j < 5$ となる j ($1 \leq j \leq R$) が存在する確率 α で表される。われわれはこの α を小さく、例えば1%未満に押さえ込みたい。つまり、 $\alpha < 1\%$ となるための条件を平均値 $\lambda_j = Np_j$ の言葉で記述すること、それに加えて、なるべく見やすい十分条件の例をいくつか見出すこと、これらがわれわれの主たる問題であり、次節で論じられる。以下セルの並べ方は、 λ_j の小さい順 ($\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_R, \sum_{j=1}^R \lambda_j = N$) とする。

危険率 α の評価式として、第一に思いつく包除原理([3])の適用により、次式を得る。

$$(1-2) \quad \sum_{j=1}^R P(x_j < 5) - \sum_{j_1 < j_2} P(x_{j_1} < 5, x_{j_2} < 5) \leq \alpha \leq \sum_{j=1}^R P(x_j < 5)$$

(1-2)の右側の不等式は、 $\bigcup_{j=1}^{2m} \{x_j < 5\} = \bigcap_{i=1}^m (\{x_{2i-1} < 5\} \cup \{x_{2i} < 5\})$ と書き直す工夫により、改良される

($m = [R/2]$; $[x]$ は小数点以下切り捨てを示すガウス記号)。

$$(1-3) \quad \alpha \leq \sum_{j=1}^R P(x_j < 5) - \sum_{i=1}^m P(x_{2i-1} < 5, x_{2i} < 5)$$

一般的に、 λ に対応するセルの人数を x , λ' に対応するセルの人数を x' と書いて、2つの量 $P(x < 5) = \beta(\lambda)$, $P(x < 5, x' < 5) = \beta(\lambda, \lambda')$ を計算できれば、 α のよい見積りに到達する。まず、単独の x は2項分布に従うので、

$$(1-4) \quad \beta(\lambda) = \sum_{k=0}^4 \binom{N}{k} \left(\frac{\lambda}{N}\right)^k \left(\frac{N-\lambda}{N}\right)^{N-k} - \sum_{k=0}^4 \frac{\lambda^k}{k!} e^{-\lambda} (= \gamma(\lambda) \text{ とおく})$$

という形に、 N が大きくて λ が普通の値のとき、ポアソン分布によって近似される。また、 $\lambda \geq 6$ のとき $\beta(\lambda)$ と $\gamma(\lambda)$ は、両者とも下に凸な単調減少関数であり、不完全ガンマ関数 $\gamma(\lambda) = \frac{1}{\Gamma(5)} \int_{\lambda}^{\infty} t^4 e^{-t} dt$ の表示式 ([2]) も容易に示される。

フェラーの確率論の本 [3]、第6章によれば、次の不等式が成り立つ。

$$(1-5) \quad \left\{ \sum_{k=0}^4 \frac{\lambda^k e^{-\frac{k^2}{N-k}}}{k!} \right\} e^{-\frac{N\lambda}{N-\lambda}} < \beta(\lambda) < \left\{ \sum_{k=0}^4 \frac{\left(\lambda e^{\frac{\lambda}{N}}\right)^k}{k!} \right\} e^{-\lambda}$$

他方、 $\beta(\lambda, \lambda')$ については、3項分布に従う x と x' の相関が非常に弱いので、2重ポアソン分布によって近似される ([3])。

$$(1-6) \quad \beta(\lambda, \lambda') = \sum_{k_1=0}^4 \sum_{k_2=0}^4 \frac{N!}{k_1! k_2! (N-k_1-k_2)!} \left(\frac{\lambda}{N}\right)^{k_1} \left(\frac{\lambda'}{N}\right)^{k_2} \left(\frac{N-\lambda-\lambda'}{N}\right)^{N-k_1-k_2} - \sum_{k_1=0}^4 \sum_{k_2=0}^4 \frac{\lambda^{k_1} (\lambda')^{k_2}}{k_1! k_2!} e^{-\lambda-\lambda'} = \gamma(\lambda) \gamma(\lambda')$$

この近似についても (1-5) と類似の不等式を導くことができる。

$$(1-7) \quad \left\{ \sum_{k_1=0}^4 \sum_{k_2=0}^4 \frac{\lambda^{k_1} (\lambda')^{k_2}}{k_1! k_2!} e^{-\frac{(k_1+k_2)^2}{N-k_1-k_2}} \right\} e^{-\frac{N(\lambda+\lambda')}{N-\lambda-\lambda'}} < \beta(\lambda, \lambda') < \left\{ \sum_{k_1=0}^4 \frac{\left(\lambda e^{\frac{\lambda+\lambda'}{N}}\right)^{k_1}}{k_1!} \right\} \left\{ \sum_{k_2=0}^4 \frac{\left(\lambda' e^{\frac{\lambda+\lambda'}{N}}\right)^{k_2}}{k_2!} \right\} e^{-(\lambda+\lambda')}$$

医療情報の開示に伴って生じ得る患者の、上述した形でのプライバシー保護に関する研究課題は、今春本学医療情報部の木村通男教授から御教示を受けたものである。著者は当初、正規近似 ([4]) に基づく結果を得て、その概略を木村教授に報告した。その後 [3] を再々読しつつ、われわれの状況 ($N=150000$, $\alpha < 1\%$) においては、正規近似よりもポアソン近似の方が、格段にすぐれていることを見出し、今秋この小論にまとめ直すことにした。この場を借りて、木村教授に感謝申しあげたいと思います。

§ 2. 患者個人が特定されるリスクの確率評価

$N = 150000$ として、患者個人が特定される危険率 α を考える。この α が 1% 未満になるには、 j 番目のセルに分割される人数 x_j の平均値 $\lambda_j (j = 1, 2, \dots, R)$ は、如何なる条件を満たせばよいか。前節のポアソン近似を活用して、この問題を論じる。厳密な条件よりもむしろ、現実にある医療情報のデータに取り組む際に、指針となるべき条件を探るのがわれわれの目的である。従って、連続量 λ は [2] の如く、半整数に限って変化するものとする。(連続量 λ で考え、計算するときは註記する。) また、確率 1% の 1%、つまり 10^{-4} 未満の量は、便宜的に無視できるものとする。

まず λ の範囲は、 $12 \leq \lambda \leq 28$ に限定してもよいという事実から始める。もし $\lambda \leq 11.5$ に対応するセルが 1 つでもあれば、(1-5) の下限の値 ($\lambda = 11.5$ のとき) を計算して、 $\alpha > \beta(\lambda) \geq \beta(11.5) > 1.07361\%$ となってしまう。

(註1.) $\beta(\lambda) = 1\%$ の解は、事実 $11.604 < \lambda < 11.605$ の範囲内にある。

他方 $\lambda \geq 28.5$ に対応するセルについて、今度は (1-5) の上限の値 ($\lambda = 28.5$ のとき) を計算して、 $\beta(\lambda) \leq \beta(28.5) < 1.33391 \times 10^{-8}$ を導く。このようなセルは高々 $[150000/28.5] = 5263$ 個しか存在しないので、 α への寄与は $1.33391 \times 10^{-8} \times 5263 \div 7.02 \times 10^{-5}$ 未満、無視できる大きさである。

こうして、 λ は $12 \leq \lambda \leq 28$ の範囲にある 33 個の半整数を動くことになる。われわれの確率評価に必要な数値 $\gamma(\lambda)$ および (1-5) の上限、下限の値については、表 1 を作成し、この節の終わりに記載する。この表は $\beta(\lambda)$ の近似値として、 $\gamma(\lambda)$ がすぐれていることを教えてくれる。すなわち、上限値と $\gamma(\lambda)$ の相対誤差は、 λ とともに増加するが、 7.2×10^{-4} を超えることはない。同様に、下限値と $\gamma(\lambda)$ の相対誤差も λ とともに増加し、上限値の場合に比して、幾分大きくなるけれども、 5.3×10^{-3} 以下に留まる。最後に、 $\beta(\lambda, \lambda')$ の近似値として 2 重ポアソン分布 $\gamma(\lambda) \gamma(\lambda')$ がすぐれている点も、(1-7) によって確認することができる。

かくして、 α は $\gamma(\lambda)$ の和として評価してもよいこととなる。何故なら、 α の評価式 (1-2) (1-3) における $\beta(\lambda, \lambda')$ の寄与分は、次式によって 10^{-4} には達しないことがわかるから。

$$(2-1) \quad \sum_{i=1}^m \beta(\lambda_{2i-1}, \lambda_{2i}) < \sum_{j_1 < j_2} \beta(\lambda_{j_1}, \lambda_{j_2}) \sim \sum_{j_1 < j_2} \gamma(\lambda_{j_1}) \gamma(\lambda_{j_2}) \leq \frac{1}{2} \left\{ \sum_j \gamma(\lambda_j) \right\}^2 \sim \alpha^2 / 2 = 0.5 \times 10^{-4}$$

今や準備完了。28 を超える λ_j は α の見積りから除外してもよいので、

$$(2-2) \quad \alpha \sim \sum_{j=1}^r \gamma(\lambda_j) < 1\%, \quad 12 \leq \lambda_1 \leq \dots \leq \lambda_r \leq 28, \quad \sum_{j=1}^r \lambda_j \leq 150000$$

を満たす数の組 $\{\lambda_j\}$ が、われわれの問題の解に他ならない。ここで r は、28 を越える λ_j を除いた残りの個数を示し、不定である。 r の値を 1 つ固定したとき、(2-2) の解全体を D_r と書けば、 D_r は凸領域をなす。つまり、 $\{\lambda_j\}$ と $\{\lambda'_j\}$ が D_r の点ならば、これらを結ぶ線分 $\{t\lambda_j + (1-t)\lambda'_j\} (0 < t < 1)$ も D_r に含まれる。従って、 D_r の境界上に位置する端点が特に重要である。まず 1 つの λ_{j_0} に着目し、

残りの $\lambda_j (j \neq j_0)$ は固定する。そして λ_{j_0} の値を小さくすれば1%の水準を越え、大きくすれば D_r 内に入るとなるような点である。(端点の凸結合で書かれる $\{\lambda_j\}$ はすべて、(2-2)の解である。)

以下このような点をいくつか構成しよう。その際 (λ が0.5の間隔で変化するため)、 α の値が1%ちょうどから $\pm 0.001\%$ の範囲内におさまるならば、限界値として十分であると考え。すなわち、次の条件で規定される $\{\lambda_j\}$ に着目しよう。

$$(2-3) \quad \left| \sum_{j=1}^r \gamma(\lambda_j) - 1\% \right| < 0.001\%, \quad 12 \leq \lambda_1 \leq \dots \leq \lambda_r \leq 28, \quad \sum_{j=1}^r \lambda_j \leq 150000$$

最小値 $\lambda_1 \geq 23$ の場合から始める。 $\lambda_1 = \lambda_2 = \dots = \lambda_{6972} = 23$ とすれば、 $\alpha = 10^{-2} - 4.8 \times 10^{-7}$ を得るが、 $23 \times 6972 > 150000$ で(2-3)の解にはなり得ない。

(註2.) $\lambda_1 = \dots = \lambda_{6544} = 22.924$ とすれば、 $\alpha = 10^{-2} - 1.8 \times 10^{-7}$ を得るが、 $22.924 \times 6544 = 150014$ で(2-3)の最後の式が不成立。($\lambda_1 \geq 22.924$ の場合、 $\sum_{j=1}^r \lambda_j \leq 150000$ から $\alpha < 1\%$ が従う。)

例1. ($\lambda_1 = 22$ の場合) $\lambda_1 = \dots = \lambda_{3038} = 22$ とすれば、 $\alpha = 10^{-2} + 9.3 \times 10^{-6}$ で(2-3)が成立。

例2. ($\lambda_1 = 20$ の場合) $\lambda_1 = \dots = \lambda_{590} = 20$ とすれば、 $\alpha = 10^{-2} - 2.6 \times 10^{-6}$ で(2-3)が成立。

例3. ($\lambda_1 = 18$ の場合) $\lambda_1 = \dots = \lambda_{118} = 18$, $\lambda_{119} = 18.5$, $\lambda_{120} = 20.5$ とすれば、 $\alpha = 10^{-2} + 7.0 \times 10^{-7}$ で(2-3)が成立。

例4. ($\lambda_1 = 16$ の場合) $\lambda_1 = \dots = \lambda_{24} = 16$, $\lambda_{25} = 16.5$, $\lambda_{26} = 17.5$ とすれば、 $\alpha = 10^{-2} + 7.8 \times 10^{-6}$ で(2-3)が成立。

例5. ($\lambda_1 = 14$ の場合) $\lambda_1 = \dots = \lambda_5 = 14$, $\lambda_6 = 15$, $\lambda_7 = 17.5$ とすれば、 $\alpha = 10^{-2} + 7.8 \times 10^{-6}$ で(2-3)が成立。

例6. ($\lambda_1 = 12$ の場合) $\lambda_2 = 14$, $\lambda_3 = 15.5$ とすれば、 $\alpha = 10^{-2} - 7.6 \times 10^{-6}$ で(2-3)が成立。

次に例6のように、 λ_1 と λ_2 の間に開きがあるケースを2、3取りあげる。

例7. ($\lambda_1 = 12$, $\lambda_2 = 24$ の場合) $\lambda_2 = \dots = \lambda_{3867} = 24$ とすれば、 $\alpha = 10^{-2} + 5.2 \times 10^{-7}$ で(2-3)が成立。

(註3.) $\lambda_1 = 12$, $\lambda_2 \geq 24.545$ の場合、註2と同じ推論により、 $\sum_{j=1}^r \lambda_j \leq 150000$ から $\alpha < 1\%$ が従う。

例8. ($\lambda_1 = 12$, $\lambda_2 = 20$ の場合) $\lambda_2 = \dots = \lambda_{143} = 20$ とすれば、 $\alpha = 10^{-2} + 6.5 \times 10^{-6}$ で(2-3)が成立。

例9. ($\lambda_1 = 12$, $\lambda_2 = 16$ の場合) $\lambda_2 = \dots = \lambda_7 = 16$ とすれば、 $\alpha = 10^{-2} + 3.0 \times 10^{-6}$ で(2-3)が成立。

以下 $\lambda_1 = 14$ で $\lambda_2 = 22$ の場合という風に、(2-3)を満たす $\{\lambda_j\}$ の例を沢山構成することができるが、一端筆を止めて、点から線へ目を転じよう。すなわち、 $\{\lambda_j\}$ の中で、2つのパラメータを取りあげ、連続的に変化させる(残りの λ_j はすべて一定に保つ)。そして1つの境界線を描くような例を構成して、この節を閉じることにしたい。

最小値 $\lambda_1 = 13$ を固定すると、 $1 - \gamma(\lambda_1) = 0.625981\%$ である。表1に加えて、 $13 < \lambda < 13.5$ の範囲内で、間隔を0.05に細分して λ を変化させ、 $\gamma(\lambda)$ の値を計算した結果(表2)をまず用意する。その後、 $13 \leq \lambda_2 \leq 13.25$ の範囲内で λ_2 の値を0.05の間隔で動かし、 $|0.625981\% - \gamma(\lambda_2) - \gamma(\lambda_3)| < 0.001\%$ を満たす実数解 λ_3 ($\lambda_2 \leq \lambda_3 \leq 13.25$) を順次求めて行く。そのとき表2だけでなく、次の5つの計算結

果(%表示)も必要となる： $\gamma(13.55) = 0.251112$, $\gamma(13.475) = 0.265219$, $\gamma(13.41) = 0.278062$, $\gamma(13.295) = 0.302269$, $\gamma(13.245) = 0.313416$ 。こうして、次の結果に到達する。

例10. ($\lambda_1 = 13$ の場合)6個の点(λ_2, λ_3) = (13, 13.55), (13.05, 13.475), (13.1, 13.41), (13.15, 13.35), (13.2, 13.295), (13.245, 13.25) を結ぶ折れ線は、下に凸で D_3 の境界線 $\gamma(\lambda_2) + \gamma(\lambda_3) = 0.625981\%$ を近似する。なお、パラメータ λ_2 の間隔0.05をもっと小さく(例えば0.01に)すれば、一段と精密な近似曲線に至る。

表1. (ポアソン近似) $N=150000$ とし、 $1 \leq a < 10$ の数 a と指数 E によって、 $a \times 10^E$ の形に表示する。

λ	$\gamma(\lambda)$	(1-5)の上限	(1-5)の下限	E
12	7.60039	7.60259	7.59241	3
12.5	5.34551	5.34713	5.33945	3
13	3.74019	3.74137	3.73563	3
13.5	2.60434	2.60520	2.60094	3
14	1.80525	1.80587	1.80272	3
14.5	1.24604	1.24646	1.24418	3
15	8.56641	8.56958	8.55277	4
15.5	5.86725	5.86951	5.85731	4
16	4.00438	4.00597	3.99717	4
16.5	2.72386	2.72498	2.71866	4
17	1.84698	1.84776	1.84325	4
17.5	1.24865	1.24920	1.24599	4
18	8.41761	8.42140	8.39864	5
18.5	5.65935	5.66197	5.64590	5
19	3.79517	3.79698	3.78568	5
19.5	2.53885	2.54010	2.53218	5
20	1.69447	1.69533	1.68980	5
20.5	1.12842	1.12900	1.12515	5
21	7.49868	7.50266	7.47593	6
21.5	4.97303	4.97573	4.95724	6
22	3.29167	3.29350	3.28074	6
22.5	2.17473	2.17597	2.16719	6
23	1.43424	1.43508	1.42905	6
23.5	9.44272	9.44836	9.40708	7
24	6.20670	6.21049	6.18230	7
24.5	4.07324	4.07579	4.05657	7
25	2.66908	2.67079	2.65772	7
25.5	1.74643	1.74757	1.73870	7
26	1.14112	1.14187	1.13587	7
26.5	7.44595	7.45100	7.41043	8
27	4.85226	4.85562	4.82825	8
27.5	3.15807	3.16029	3.14187	8
28	2.05291	2.05438	2.04200	8

表2. $13.05 \leq \lambda \leq 13.45$ における $\gamma(\lambda)$ の値(%表示)

λ	13.05	13.1	13.15	13.2	13.25	13.3	13.35	13.4	13.45
$\gamma(\lambda)$	0.360799	0.348031	0.335698	0.323787	0.312284	0.301175	0.290448	0.280090	0.270090

謝辞

1. 今後、病院に蓄積された臨床データに取り組み、リスク評価の実証的研究を継続するよう、レフェリーから励ましの言葉をかけられる。そのような研究に伴う問題点と指針を合わせて、教示して下さいたことに対し、深甚なる感謝を申しあげます。
2. 資料の整理とともに、原稿の清書をお願いした鴨藤江利子さんに、厚く御礼申しあげます。

参考文献

- [1] 井上潔司, 安芸重雄: 一般化された誕生日問題, 日本数学会年会(於東京大学), 統計数学分科会講演アブストラクト: 85-86, 2009年3月.
- [2] 春日屋伸晶(編): 実用数表大系15 ガンマ函数表・ベータ函数表. 技報堂, 1972.
- [3] W. Feller: An Introduction to Probability Theory and Its Applications, Vol. I (3rd ed.). New York: John Wiley & Sons, 1968.
- [4] 山内二郎(編): 統計数値表 JSA-1972. 日本規格協会, 1972.