

連載

## EBNのための統計の読み方⑫

# EBNのための研究計画

## 統計学的推論：推定と検定

林 邦彦・高木廣文

# EBN のための研究計画

## 統計学的推論：推定と検定

林 邦彦 HAYASHI Kunihiko

群馬大学医学部保健学科

高木廣文 TAKAGI Hirofumi

新潟大学医学部保健学科

「ハリー、君は何かを学ばなかつたかね。我々の行動の因果というものは、常に複雑で、多様なものじや。だから、未来を予測するということは、まさに非常に難しいことになるのじやよ」  
(ダンブルドア校長)<sup>1)</sup>

### 統計学的推論

前号では、調査や実験などで収集したデータについて、その「要約」や「データ構造」の記述を目的とした解析法の話であった。今号では、単に手元にあるデータの様子を把握するだけではなく、そのデータを利用して、背後にある自然法則を推論し、母集団や未来における予測に役立てることを考えてみる。

### ◆ 「推定」と「検定」

自然法則を推論する統計学の代表的な方式に、「推定 (estimation)」と「検定 (test)」がある。また、自然法則をつかさどるものや、集団の真の特性を示すものを「母数 (parameter)」という。日常会話で、「ヨーグルトが6割の人で効くっていっていたけれど、母数が5人の調査じゃ信用できないよね」などというように、データ数や分数の分母 (denominator) のことを「母数」という人がいるが、統計学の用語としては誤用といえる。自然法則をつかさどる母数の値は、残念ながら直接に測定したり観察したりする

ことができない。そこで、母数の値を手元のデータから何とか推し測ろうとするのが「推定」である。

一方、相反する2つの仮説を用意して、手元のデータから正しそうな仮説を選択するという、二者択一的な方式で自然法則を推論するのが「検定」である。

### ◆ 「母集団」と「標本」

ここで、統計学の教科書でまず登場する、「母集団 (population)」と「標本 (sample)」について考えてみよう。

総勢520人の看護師が勤務している病院看護部で、看護師が使う手洗い用洗浄液として製品Aと製品Bのどちらを使うか、意見を聞いて決めようと思う。全員聞くのは困難なので、全体の1/20の26人に意見を聞くことにした。

この場合、「母集団」は520人の看護師全員と明示的に表すことができる。意見を聞くのに選ばれた26人が「標本」である。この標本のデータから統計学的推論を用いて、全員の様子を推測し、一般化する。そのため、標本26人は看護部全員の様子を代表するように選ぶ。全員から無作為に抽出したり、病棟ごとの看護師数に応じて代表者を選ぶなどの工夫をする。世論調査や一般住民対象の有病調査の場合などには、選挙人名簿や住民台帳から母集団を明確に定義して、層化多段階無作為抽出法などで標本を選び出す。

しかし、医学分野の多くの研究では、母集団から抽出し、母集団の特徴を保持した「標本」に基づくデータとは単純にはいかない。たとえば、米国のH大学病院に入院した26人の肺癌患者を対象に、治療法Aと治療法Bのどちらが効果で優れているかを調べ、統計学的推論を用いた研究があるとしよう。標本数は前例と同じ26人だが、この研究の母集団は一体、何だろう。対象になった入院患者のほかに、H大学病院に外来通院中の肺癌患者や、昨日新たに肺癌と診断された患者は母集団に含まれるのか。また、米国内の他病院や日本の病院に入院している肺癌患者などは、どうであろう。前述の520人の母集団の例とは違い、大変にあいまいである。母集団が明示できないので、標本を無作為に抽出することもできない。統計学的推論さえ使えば、すぐさま何にでも推論結果を一般化できるわけではないことに注意してほしい。

### ◆統計学的推論における一般化

医学研究では、推論の一般化について、次のように段階的に集団を分類して考えることが多い<sup>2)</sup>。「白人の肺癌患者」といった研究結果を適用できると考えるいくつかの集団群を「想定母集団群（target populations：標的集団）」という。想定母集団群のなかで「2003年にH大学病院に入院した肺癌患者」といった実際のデータ源となる集団を「データ源集団（source population）」という。データ源集団のなかの「PSが2で、80歳以下の転移がない肺癌患者」といった研究対象の基準を満たす集団を「対象適格集団（eligible population）」、そして対象適格集団のなかでも「研究参加に同意し実際にデータを提供する患者」が「研究参加者集団（participant population）」となる。もう、population（集団）だらけである。結果を一般化しようとする想定母集団群と、実際にデータが収集される研究参加者集団との間の各段階で、特徴の保持をゆがめる可能性があることがわかるであろう。前述の母集団と標本のような單

純な関係ではない。しかし、統計学的推論ではそれを前提としている。

## 八 推定の考え方

### ◆分布の性状

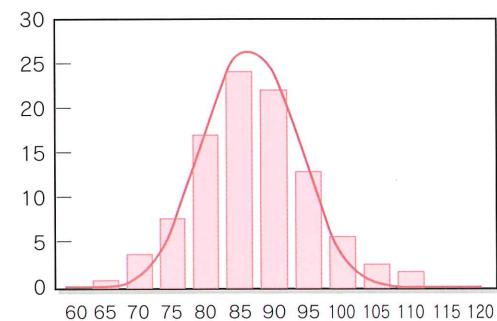
さて、実際に手元にあるデータを使って、「母数」を推定してみてみよう。

糖尿病や過度の肥満の人を除いた健康な30歳代女性100人の空腹時血糖値のデータを利用して、30歳代の一般健康女性集団（母集団）では空腹時血糖値はどんな値となるか推測する。①に、標本100人の空腹時血糖値の度数分布を示す。

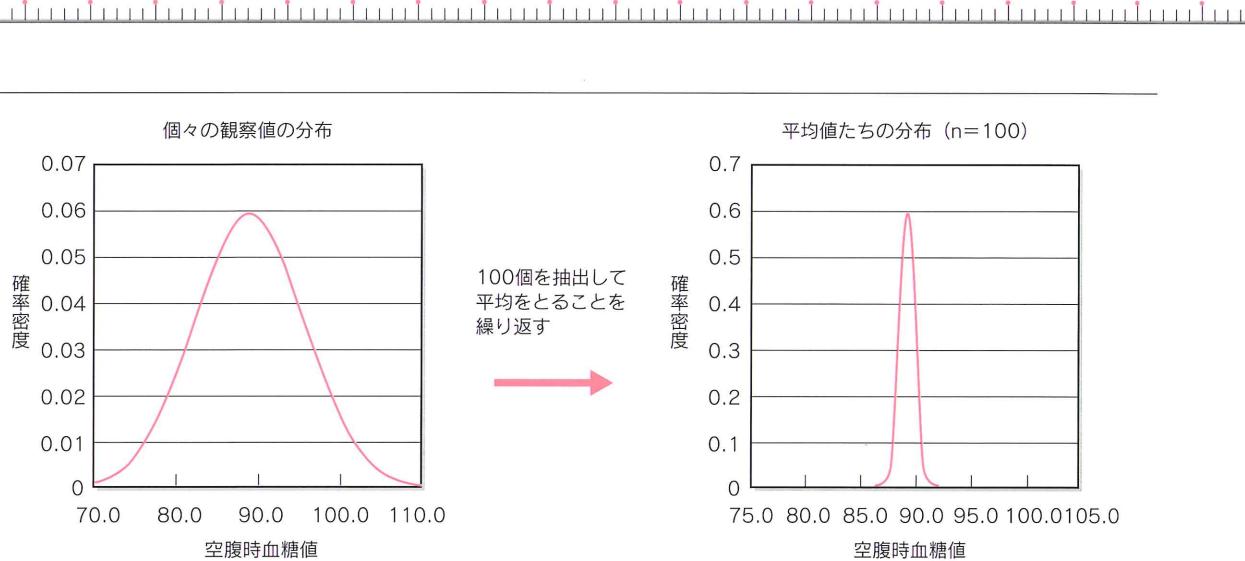
「標本データ」の平均値 $\bar{x}$ は88.97（標本平均）、不偏分散 $s^2$ は $(6.73)^2$ である。これらを使って「母集団」での値の分布の様子を推測するには、母集団での平均である「母平均 $\mu$ 」を「標本平均 $\bar{x}$ 」で、母集団での「母分散 $\sigma^2$ 」を「不偏分散 $s^2$ 」で推定する。 $\sigma$ が標準偏差である（前号を参照）。もし、母集団で各人の空腹時血糖値が、①の曲線のように左右対称鈎鐘型の正規分布（normal distribution）としてばらつくならば、その様子は

確率密度関数

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



①空腹時血糖値の度数分布と正規分布近似



## ② データの分布とその平均値の分布

と表すことができる。ただし、 $\pi$ は円周率を示す。これは、そもそも天文学での星の観測の誤差を説明するために、Gaussが提案したので、ガウス分布ともよばれる<sup>3)</sup>。位置を決める「母平均 $\mu$ 」と、ばらつき度合いを決める「母分散 $\sigma^2$ 」(もしくは標準偏差 $\sigma$ )の2つの未知母数を推定できれば、この分布の性状を決めることができる。

実際のデータから推定した分布が②の左図である。横軸に $x$ 、縦軸に $f(x)$ をとってある。曲線下の全面積はちょうど1となる。また、 $\mu \pm \sigma (= 88.97 \pm 6.73 : 82.24 \sim 95.70)$  の間の曲線下面積は0.682、 $\mu \pm 1.96 \sigma (= 88.97 \pm 13.19 : 75.78 \sim 102.16)$  の間の曲線下面積は0.950となる。つまり、平均値から標準偏差の大きさを上下にとった区間に全体の68.2%が、標準偏差の1.96倍の大きさを上下にとれば全体の95%が分布する。

### ◆ 標準誤差と信頼区間

次に、この母集団の分布から、今回の標本数である100個のデータを選んでは平均値を計算することを繰り返し行なったことを考えてみよう。毎回算出する平均値は、100個のデータの選び方次第で少しずつ異なるだろう。平均値たちを無数に算出できたとして、

その平均値たちはどのように分布するかを、正規分布を用いて推測したものが②の右図である。その分布の中心の位置は、やはり母平均の推定値である標本平均88.97にくるだろう。

左図と異なる点は、ばらつき度合いである。平均値たちの分散は、母分散 $\sigma^2$ を標本数 $n$ で割った $\frac{\sigma^2}{n}$ である。その平方根をとった $\frac{\sigma}{\sqrt{n}}$ を「標準誤差(SEM: standard error of mean)」とよぶ。

100人の空腹時血糖データからは、平均値の標準誤差は $6.73/10 = 0.673$ である。標準誤差の式からわかるように、同じ母集団から抽出した標本による平均値でも、その標準誤差の値は標本数によって異なることがわかる。標本数が2倍になると、標準誤差は $\frac{1}{\sqrt{2}}$ 倍になり小さくなる。標本数が大きなデータ・セットから推定したほうが平均値たちのばらつきは小さくなり、②の右図での分布の山はより尖ったものになる。つまり、標本数が大きいほど算出した平均値たちの揺るぎが少ないと考えていいだろう。この平均値たちの揺るぎの程度、別の表現をとれば、算出した平均値の確信の程度を示す方法に、「信頼区間(CI: confidence interval)」がある。

$(100 - \alpha)\%$ 信頼区間は、その区間に母数の真値が $(100 - \alpha)\%$ の確率で含んでいることを示し、

$$\bar{x} - Z(\alpha\%) \cdot \text{SEM} \leq \mu \leq \bar{x} + Z(\alpha\%) \cdot \text{SEM}$$

によって求められる。 $(100 - \alpha)\%$ は「信頼係数」といい、95%がよく用いられる。 $Z(\alpha\%)$ は標準正規分布<sup>注1)</sup>の両裾の面積が $\alpha\%$ となる点の値を表し、 $Z(10\%) = 1.645$ 、 $Z(5\%) = 1.960$ 、 $Z(1\%) = 2.576$ である。上記の100人の空腹時血糖値から推定した母平均の95%信頼区間は、

$$88.97 - 1.960 \times 0.673 < \mu < 88.97 + 1.960 \times 0.673$$

となり、87.65～90.29の間に95%の確率で母平均があると推定される。確信の程度が低くてよければ90%信頼区間を、高くしたければ99%信頼区間を計算すればよい。

このように、母平均の位置について、標本での平均値という一つの値でぴたりと推定する方法を「点推定」といい、95%信頼区間のような値の幅で推定する方法を「区間推定」という。

この信頼区間は、「母比率」の推定にも利用できる<sup>4)</sup>。先ほどの30歳代女性100人のうち、現在喫煙をしている人は12人であった。この標本での喫煙割合 $p$ は $12/100 = 0.12(12\%)$ であるが、この値は母集団での喫煙の「母比率 $\pi$ 」の推定値となる。この「標本比率 $p$ 」の分散は $p(1-p)$ 、標準偏差 $\text{SE}(p)$ は $\sqrt{\frac{p(1-p)}{n}}$ となる。そこで、母平均のときと同じように、母比率 $\pi$ の95%信頼区間は、

$$p - Z(\alpha\%) \cdot \text{SE}(p) \leq \pi \leq p + Z(\alpha\%) \cdot \text{SE}(p)$$

によって求められ、

$$0.12 - 1.960 \times 0.0325 < \pi < 0.12 + 1.960 \times 0.0325$$

と、母集団の30歳代女性での喫煙母比率は0.0563～0.1837（つまり、5.63%～18.4%の間）にあると、95%の確率で推定できる。

注1) 標準正規分布：正規分布のなかでも平均0、分散1の分布を「標準正規分布」（standard normal distribution）とよぶ。統計学の教科書の多くに、この標準正規分布のZ値と $\alpha\%$ が、数表として載っている。

## 八 検定の考え方

### ◆帰無仮説と対立仮説

次に、もう一つの統計学的推論の方式である「検定」に話を移そう。

前述のH大学病院での肺癌治療法の研究の例で、治療法A・Bの違いが改善・悪化という病状の変化に関係があるか推論することを考えてみよう。このような研究での作業仮説は、今まで標準的に用いられてきた治療法Aに比べて、新治療法Bが優れていることを示すことなのかもしれない。しかし、「統計学的検定」では相反する2つの仮説、「帰無仮説 $H_0$ （null hypothesis）」と「対立仮説 $H_1$ （alternative hypothesis）」を用意する。2つの仮説を合わせると、すべての状況が記述されているように設定する。肺癌治療法の例であれば、

帰無仮説  $H_0$ ：治療法と病状は独立している

（治療法間で病状は同じ）

対立仮説  $H_1$ ：治療法と病状は従属している

（治療法間で病状が異なる）

とすれば、この2仮説以外には治療法と病状の関係を示す状況はなくなる。であれば、この2仮説のいずれか正しいほうを一定ルールのもとで選ぶことができれば、真実が推論できる。

仮説選択のルールは、あたかも二重否定のように回りくどい。まず、自分の主張したい仮説の逆である「帰無仮説」の立場をとってみる。帰無仮説 $H_0$ が正しいと仮定して、今回得られたデータはどの程度起きやすいのかを、「検定統計量」（t検定でいえばt値が検定統計量）からみてみる。算出した検定統計量が事前に設定していた限界値を超える、もしくは検定統計量から算出したp値が小さくて、今回のデータが示す状況は帰無仮説 $H_0$ のもとでは極めて起き難



	病状改善(人)	病状悪化(人)	計(人)
治療法A	3	10	13
治療法B	9	4	13
計	12	14	26

### ③米国のH大学病院での肺癌治療の成績（仮想データ）

いこと（たとえば、起こる確率が5%未満）と考えられた。しかし、それは大変まれなことが起きたと考えるよりは、前提としていた帰無仮説 $H_0$ が誤っていたとして捨て去る（棄却：rejectという）ほうが合理的であろう。帰無仮説 $H_0$ を捨てると生き残った仮説は今や対立仮説 $H_1$ のみとなって、自分の言いたかった「治療法間で病状が異なる」という対立仮説 $H_1$ が正しいことが検証される、という筋書きである。ややこしい筋書きだが、まず自分の意見と反対の立場に立って考えてみることが、統計的検定だけではなく社会一般で重要なようである。

この論法を、肺癌治療の例（③）で $\chi^2$ 検定<sup>注2)</sup>を使って追ってみる。まず、帰無仮説が正しいとして、今回の状況はどの程度起きやすいのか調べる。Yatesの補正した $\chi^2$ 値は、

$$\chi_0^2 = \frac{26(|3 \times 4 - 10 \times 9| - \frac{26}{2})^2}{13 \times 13 \times 12 \times 14} = 3.869$$

となる。この値は、事前に設定した自由度1の有意水準5%である $\chi^2$ 値の限界値 $\chi^2(1, 0.05) = 3.841$ を超えていた。また、 $\chi_0^2$ 値から $p$ 値を計算した場合でも $p = 0.0492$ となり、これは帰無仮説の下では5%未満でしか起きない状況であることがわかる。そこで、前提としていた帰無仮説 $H_0$ が誤っていたとして棄却する。すると、対立仮説 $H_1$ が残り、治療法間で病状は統計学的に有意に異なると推論することになる。

		検定の結果	
		$H_0$ を棄却できず	$H_0$ を棄却（ $H_1$ を採用）
真実	帰無仮説 $H_0$	$1 - \alpha$	$\alpha$
	対立仮説 $H_1$	$\beta$	$1 - \beta$

$\alpha$ ：第1種の過誤、 $\beta$ ：第2種の過誤

### ④ 検定における2種類の過誤の確率

これで、H大学病院では、胸を張って治療法Bを勧めるようになるかもしれない。

しかし、この論法の弱点は、帰無仮説が棄却できないときに表われる。二重否定的な論法になつてないので、否定できなかつた（帰無仮説 $H_0$ を棄却できなかつた）場合には、論理が途中で止まり、態度保留となつてしまう。帰無仮説 $H_0$ が棄却できないからといって、積極的に帰無仮説 $H_0$ が正しいとはいえないものである。

### ◆ 検定で起こり得る2種類の過誤

実際、帰無仮説 $H_0$ が棄却できないときには、眞実 $H_0$ が正しい場合と、実は対立仮説 $H_1$ が正しいのにそれが検出できない場合がある。④を参照してほしい。

検定には、2種類の推論の誤りが起こり得る。第1種の過誤（type-I error）の確率 $\alpha$ は、帰無仮説 $H_0$ が眞実正しいにもかかわらず、誤って棄却してしまう確率である。検定で仮説選択の基準とする5%といった有意水準（significant level）は、この確率の大きさのことである。このように、第1種の過誤の大きさは通常は5%未満におさえられていることになる。

検定でのもう一方の推論の誤りが、第2種の過誤（type-II error）である。第2種の過誤の確率 $\beta$ は、

注2)  $\chi^2$ 検定：2×2分割表において独立性の検定として用いられる。n数が小さな場合には、フィッシャーの直接確率法が用いられるが、n数がある程度以上あれば、Yatesにより考案された連続修正（ $\chi^2$ の式の分子の $-\frac{n}{2}$ の部分）した $\chi^2$ 検定が用いられる。

また、1×m分割表での自由度は $(l-1) \times (m-1)$ であり、2×2分割表では自由度は1となる。

対立仮説  $H_1$  が真実正しいにもかかわらず、帰無仮説を棄却できない確率となる。あまりにも標本数が小さな研究や、測定や観察が不正確な研究では、この見逃しの確率  $\beta$  が大きくなってしまう。真実は対立仮説  $H_1$  が正しいときに、検定で帰無仮説  $H_0$  を棄却し、対立仮説  $H_1$  が採択できる確率  $1 - \beta$  を「統計学的検出力 (statistical power)」という。研究計画では、この検出力が 80~90% 程度になるように設計される。

## ◆両側検定と片側検定

また、検定には「両側検定」と「片側検定」がある。前述の癌治療法の例では、研究の作業仮説は、治療法 A に比べて新治療法 B が優れていることであった。しかし、設定した対立仮説は「治療法間で病状が異なる」である。「治療法間で病状が異なる」には、「治療法 B のほうが治療法 A より病状がよい」と「治療法 B のほうが治療法 A より病状が悪い」の 2 つのケースがある。両方のケースを含んで、対立仮説を「治療法間で病状が異なる」とすれば、それは「両側検定 (both sided test)」をすることになる。新治療 B のほうがずいぶんよいだろうと思って実施した研究でも、実は今までの治療法 A のほうがよいということも起こり得る。多くの場合、両側検定が使われる。しかし、理論的に片方のケースしか想定

できない場合には、対立仮説を「治療法 B のほうが治療法 A より病状がよい」、もしくは「治療法 B のほうが治療法 A より病状が悪い」のいずれかにして「片側検定 (one sided test)」を行うことになる。

## おわりに

統計学的推論として、「推定」と「検定」の考え方を解説した。医学論文や学会発表の図表で、検定結果を有意水準 5% で有意であった個所に \*, 1% で有意を \*\*, 0.1% で有意を \*\*\* と表示し、「差は小さくても、統計学的に有意」とか、「星 3 つです。きわめて有意です」など、まるで食通のレストラン情報のように使われることもある。しかし、看護や医学の研究では、有意か否か二者択一方式の推論より、差はどの程度といった大きさに関する母数の推定値のほうが重要な情報となることが多い。

また、検定では、帰無仮説が棄却できて初めて有効な仮説選択ができる、帰無仮説が棄却できないときは選択保留とならざるを得ない。同じであることを積極的に検定で示すには、「同等性の検定」や「非劣性の検定」とよばれる特殊な仮説構造が必要となる。ぜひ「推定」と「検定」という 2 つの異なる統計学的推論の長短所を理解して、上手に使い分けてほしい。

## 文献

- 1) ローリング J.K.: Harry Potter and the prisoner of Azkaban. 松岡祐子訳: ハリー・ポッターとアズカバンの囚人. 静山社; 2001.
- 2) 林邦彦: コホート研究の読み方・チェックリスト. EBM ジャーナル 2001; 2(5): 60-63.
- 3) Walker AM: Observation and inference—An introduction to the methods of epidemiology. Epidemiology Resources Inc., MA; 1991. 丸井英二, 中井里史, 林邦彦訳: 疫学研究の考え方・進め方—観察から推測へ. 新興医学出版; 1996.
- 4) 柳井晴夫, 高木廣文編: 新版看護学生書基礎科目 統計学. メヂカルフレンド社; 1995.