

## 検定

---

前回説明した区間推定では、母平均について、1つの数で推定せずにある程度の余裕をもって推定しました。しかし、実際には母平均などのパラメータについて「ある値であるかないか」「ある値より大きい(小さい)か否か」などの判断が必要なことがしばしばあります。検定(仮説検定)とは、統計的推測の考え方を使って、母集団に対する判断を行うものです。その基本的考え方は、

めったに起きないことは、いま現実には起きていない  
(いま起きていることは、ありふれたできごとにはちがいない)

というものです。

いま、あるくじが50%の確率で当たるとします。このくじを10本ひくと、全部はずれでした。このとき、「50%の確率で当たるくじが10本続けてはずれることなど、めったにあるはずがない」と考えるのは自然なことです。そうすると、最初の仮定は間違いで「このくじの当選確率は50%よりも小さい」と考えるほうが妥当です。このような素朴な考え方を統計学を用いて述べたのが仮説検定です。

上のくじびきの問題に答えるには、第5回で説明した2項分布モデルを用います。今回は、2項分布にもとづいて検定の考え方を説明します。

---

### 2項分布と検定

では、冒頭でふれた状況をもう一度考えてみましょう。

あるくじは、50%の確率で当たるとします。このくじを10本ひくと、全部はずれでした。このくじは、本当に50%の確率で当たるのでしょうか。当たる確率はもっと小さいのでしょうか？

このときは、「50%の確率で当たるはずのくじが10本続けてはずれることなど、めったにあるはずがない」と考えるのは自然なことです。そうすると、最初の「50%の確率で当たる」は間違いで、「このくじの当選確率は50%よりも小さい」と考えるほうが妥当です。では、次の状況はどうでしょうか。

あるくじは、50%の確率で当たるとします。このくじを50本ひくと、17本当たって33本はずれでした。このくじは、本当に50%の確率で当たるのでしょうか。当たる確率はもっと小さいのでしょうか？

今度は少々微妙な状況です。前の状況と同じように「当たる確率は50%よりも小さい」のかもしれませんが、もしかしたら「当たる確率は50%」というのが本当で、単に運が悪かったのかもしれない。

どちらが正しいのかは、わかりません。しかし、当選確率が50%か、それともそれより小さいのか、どちらの言い分がよりもっともかを、統計学を使って判断するのが仮説検定(あるいは単に検定)という手法です。

では、仮に「当選確率は50%である」ということを認めるとしましょう。さきほど、「当選確率50%のとき、10回続けてはずれることはめったにない」ということを述べましたが、では「当選確率が50%のとき、50回ひいて17回しか当たらない」ことは、どのくらい「めったにない」ことなのでしょう。それを計算することを考えてみましょう。

計算しなければならない確率は

ここで、「そんなことはめったにない」の「そんなこと」とは、正確にはどういう意味でしょうか。それは、「50本ひいてちょうど17本だけ当たる」ことではありません。「50本ひいて、16本でも18本でもなく、ちょうど17本当たる」ことは、たしかにめったにないでしょう。しかし、われわれはそれを問題にしているわけではありません。

仮に、50本のくじをひいて1本も当たりが出なければ、「50%の確率で当たる」というのはきわめて疑わしいでしょう。それは、「50%の確率で当たる」はずなのに50本中1本も当たらない、という確率はきわめて小さいからです。確率がより小さいはずのできごとが起きると、そもそもの「50%の確率で当たる」という前提に対する疑いは、より強くなります。

ということは、「50本ひいて17本しか当たらない」ことが「めったにない」、という言い回しは、言外に、16本しか当たらないことも、15本しか当たらないことも、14本しか当たらないことも、..., 1本も当たらないことも、当然すべて「めったにない」と考えている、ということを含んでいるはずで、17本しか当たらないことが不満な人は、16本しか当たらないことも、当然1本も当たらないことも、言うまでもなく不満なのです。

つまり、計算しなければならないのは、「不満をもつような現象」が起きる確率、すなわち「50本ひいたとき、当たり本数が17本以下である」確率です。

この確率は、第5回で説明した2項分布を使って求めることができます。しかし、この計算は「階乗」の計算が入っていて結構面倒です。それに、今回もとめるのは「当選確率50%のくじを50本ひいたとき、当たり本数が17本以下である」確率ですから、この計算を、当たる回数が17回、16回、..., 0回のすべての場合について行なわなければなりません。

ド・モアブル＝ラプラスの定理

そこで、この計算を別の方法で簡単に行なうための、ド・モアブル＝ラプラスの定理というものが知られています。ド・モアブル＝ラプラスの定理は、中心極限定理を使って2項分布を正規分布で近似し、正規分布の数表を使って計算をする方法です。

「1回あたり確率 $p$ で成功する、 $n$ 回のベルヌーイ試行」を別の視点から見てもみましょう。ちょっと変な感じですが、「1回のベルヌーイ試行で、成功する回数」を考えて、 $X$ で表します。当然、 $X$ のとりうる値は0回または1回です。1回あたり確率 $p$ で成功しますから、「1回のベルヌーイ試行で、成功する回数」 $X$ は2項分布 $B(1, p)$ にしたがいます。

一方、「1回あたり確率 $p$ で成功する $n$ 回のベルヌーイ試行で、成功する回数」を $S$ とすると、 $S$ は2項分布 $B(n, p)$ にしたがう、すなわち、成功回数が $x$ である確率は2項分布 $B(n, p)$ で計算できます。しかし、見方を変えれば、この $S$ はさっきの $X$ を $n$ 個合計したものと考えることもできます。

$n$ 個の $X$ は互いに独立ですから、それらの合計である $S$ は、 $n$ が大きいときは中心極限定理によって概ね正規分布にしたがいます。一方、 $S$ は2項分布 $B(n, p)$ にしたがうのですから、前回の説明で述べたように、 $S$ の期待値は $np$ 、分散は $np(1-p)$ です。 $S$ の分布を、「概ね正規分布」とよぼうが、「2項分

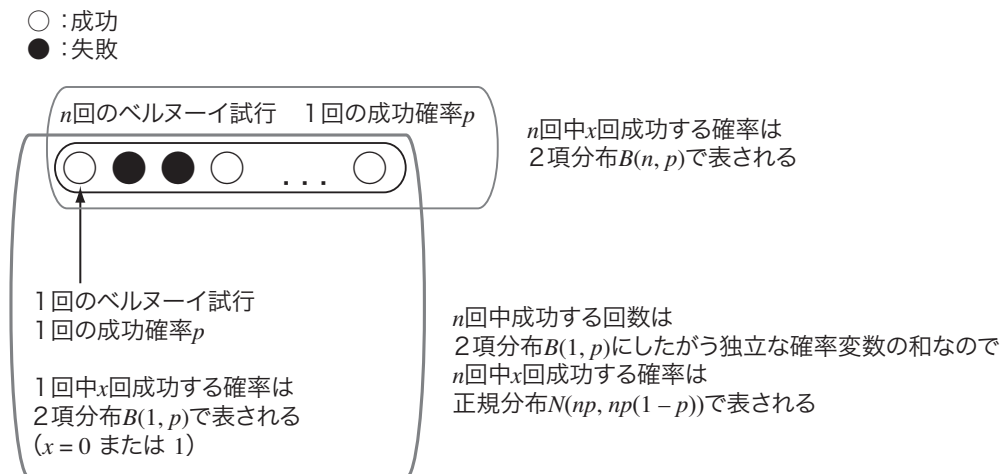


図 1: ド・モアブル＝ラプラスの定理

布」とよぼうが、同じ現象を別の名前でもよんでいるだけです。期待値や分散は同じです。したがって、 $n$  が大きいとき、 $S$  のしたがう 2 項分布  $B(n, p)$  は正規分布  $N(np, np(1-p))$  で近似できることがわかります。

以上から、「当選確率 50% のくじを 50 回ひいて、当たり回数が 17 回以下である」確率は、 $p = 0.5$ 、 $n = 50$  としたとき、正規分布  $N(np, np(1-p))$  にしたがう  $S$  が「 $S \leq 17$ 」である確率となります。 $S$  は正規分布  $N(np, np(1-p))$  にしたがうので、さらに

$$Z = \frac{S - np}{\sqrt{np(1-p)}} \quad (1)$$

とおくと  $Z$  は標準正規分布  $N(0, 1)$  にしたがいます。また、 $S = 17$  のとき

$$Z = \frac{17 - 50 \cdot 0.5}{\sqrt{50 \cdot 0.5 \cdot (1 - 0.5)}} = -2.26 \quad (2)$$

ですから、「 $S \leq 17$ 」である確率は「 $Z \leq -2.26$ 」となる確率となります。数表から、この確率は 0.0119 であることがわかります。

### $p$ 値と検定

以上のことから、「当選確率 50% としたとき、このくじを 50 回ひいて、当たり回数が 17 回以下である」確率は 0.0119 であることがわかりました。この値を  $p$  値といいます。

「当選確率 50% としたとき、このくじを 50 回ひいて、当たり回数が 17 回以下である」確率は 0.0119 です。では、「当選確率 50% としたとき、このくじを 50 回ひいて、当たり回数が 17 回以下である」ことは「めったにないこと」なののでしょうか？ それはなんとも言えません。「確率 0.0119 でしか起きない」のか、「起きる確率が 0.0119 もある」のか、それは考え方によるでしょう。

しかし、それでは話は続きません。そこで、統計学では「起きる確率がある値以下のときは、それは『めったにない』ことだ」と考えます。この値を有意水準といい、通常 5%(0.05) や 1%(0.01) が用いられます。

ここでは、有意水準を 5% としましょう。そうすると、0.0119 は 5% より小さいですから、「当選確率

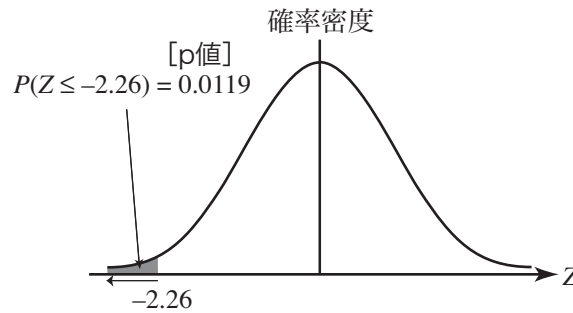


図 2:  $p$  値

50%としたとき、このくじを 50 回ひいて、当たり回数が 17 回以下である」ことは「めったにないこと」だ、と判断されます。そこで、

仮に、当選確率が 50%とする

- そうすると「50 回くじをひいて、当たり回数が 17 回以下である」という有意水準よりも小さい確率でしかおきない、つまり「めったにないはず」のことが現に起きていることになる
- 当選確率が 50%というのは間違いで、当選確率は本当はもっと小さい、と結論する。

という推論が行えます。この推論を仮説検定（あるいは単に検定）といいます。

最初の「当選確率が 50%とする」という仮定を帰無仮説といい、 $H_0: p = 0.5$  と表します。また、「当選確率が 50%というの間違いと結論する」ことを、帰無仮説を棄却するといいます。さらに、帰無仮説が棄却された結果得られる「当選確率はもっと低い」という結論を対立仮説といい、 $H_1: p < 0.5$  と表します。

#### 棄却域と検定統計量

上で、あたり回数  $S$  が 17 以下である確率は、標準正規分布にしたがう確率変数  $Z$  が  $-2.26$  以下である確率と同じで、数表からこの確率、すなわち  $p$  値は 0.0119 である、と述べました。

ところで、数表から「 $Z \leq -1.64$ 」となる確率が 5%です。ですから、 $Z$  の値が  $-1.64$  以下であれば  $p$  値が 5%以下となり、このとき帰無仮説を棄却します。すなわち、いちいち  $p$  値を計算しなくても、 $Z$  の値を求めてそれが  $-1.64$  以下であれば、有意水準が 5%のとき帰無仮説は棄却されます。この意味で、「 $Z \leq -1.64$ 」をこの問題における棄却域といい、計算の結果  $Z$  の値が棄却域に入ることを「棄却域に落ちる」といいます。また、棄却域を表すのに用いる確率変数  $Z$  を検定統計量といいます。

#### 棄却されないときは

ここまで述べてきたように、検定では、帰無仮説は「内心では」棄却されることが期待されています。目論見通り棄却されると、「対立仮説を採択する」という結論が得られるわけです<sup>1</sup>。

では、今回のくじびきの例で、有意水準を 1%にしてみましましょう。この場合、 $p$  値が 1%以下ならば帰無仮説を棄却します。ところが、今回の例では  $p$  値は 0.0119 ですから、わずかながら 1%よりも大きくなっています。したがって、目論見に反して、帰無仮説は棄却されません。

<sup>1</sup>帰無仮説を「無に帰す仮説」とよぶのはその意味です。

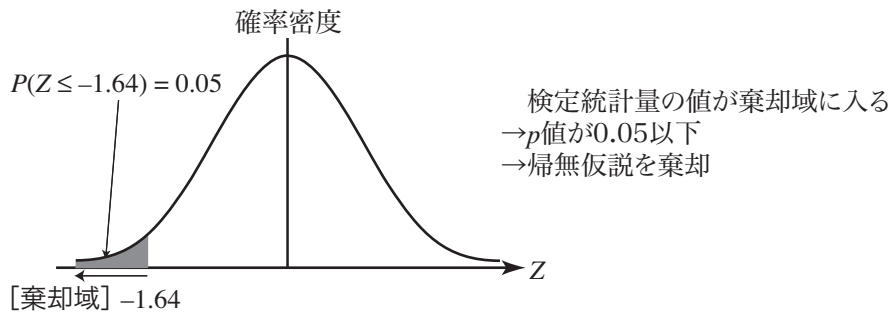


図 3: 棄却域

この場合、帰無仮説が棄却されなかったのは、「帰無仮説が正しいとき、今おきている現実得られる確率は非常に小さいとまでは言えない」です。したがって、「帰無仮説が間違っているかどうかはわからない」「対立仮説が採択できるかどうかはわからない」という結論を導かなくてはなりません。今回の例でいえば、帰無仮説が棄却されなかった場合は、「当選確率は50%より小さいとは言えない」と答えなければなりません。つまり、「目論見はずれた。当選確率は50%より小さいとまで断言する自信はない」という結論になるのです。

注意しなければならないのは、あくまで、「帰無仮説が正しいとき、今おきている現実得られる確率は非常に小さいとまでは言えない」であって、「確率が大きい」のではない、ということです。したがって、帰無仮説が棄却されなかったときに、「帰無仮説が正しい」「対立仮説は間違っている」という結論が得られるわけではありません。今回の例でも、「当選確率は50%である」などと答えてはいけません。つまり、

帰無仮説を棄却しない = 帰無仮説を採択する  
対立仮説を採択するべきかどうか断言できない

ということです。なお、「帰無仮説を棄却すべきなのに棄却しない」という誤りを第2種の誤りといいます。

### 有意水準の違い

ところで、最初のくじびきの例では、有意水準が5%のときは「帰無仮説を棄却する = 50%の確率で当たるといのは間違い」と結論され、有意水準が1%のときは「帰無仮説を棄却しない = 50%の確率で当たるといのは間違い」と言い切れない」という結論になりました。しかし、「50本中17本しか当たらなかった」という現実は同じで、有意水準は勝手に決めたのに、こんなに違った結論になってもよいのでしょうか？

これについては、「検定とはそういうものだ」ということを、よく理解しなければなりません。有意水準は、検定をする人の「大胆さ・慎重さ」の程度を表しているのです。

有意水準が大きい(5%)ときは、今起きているような現実(17本しか当たらない)が起きる確率が少々大きくても、5%以下であれば「そんなことが起きるはずがない、帰無仮説は間違っている」と結論します。はっきり物をいう態度ではありますが、帰無仮説が正しい(おじさんが正直で、偶然17本しか当た

らなかった) ときでも「間違っている」と断言してしまう可能性があります。大胆ですが、勇み足も多い、というわけです。

有意水準が小さい(1%)ときは、今起きているような現実が起きる確率が、1%以下と相当小さくないと、「わずかでもそんなことが起きる可能性があるのだから、帰無仮説は間違っているとは言い切れない」となり、結論を出しません。慎重ですが、煮え切らない態度ということになります。

検定はどんな時にするものなのか

ところで、帰無仮説が正しい(おじさんが正直で、17本しか当たらなかったのは単なる偶然だった)ときに帰無仮説を棄却してしまうという誤りを、第1種の誤りといいます。

つまり、

仮に帰無仮説が正しいとしても、有意水準**5%(1%)**の仮説検定を何度も行くと、そのうち**5%(1%)**は第1種の誤りを犯して棄却し、採択すべきでない対立仮説を採択してしまう

ことになります。

ですから、同じ現象についてなんどもデータを集めて、同じ帰無仮説について検定を繰り返す、たまに対立仮説が採択されても、直ちに「帰無仮説は間違っている」とはいえません。例えば、「血液型と性格に関係はない」という帰無仮説について何度もデータを集めて検定を行い、たまに「血液型と性格に関係がある」という結論が出ても、直ちに「やっぱり血液型と性格に関係がある」ということにはなりません。何度も検定を帰無仮説が間違っていない場合でも、たまに対立仮説が採択されるのは、むしろ自然なことです。血液型と性格の問題でいえば、ごくたまに「血液型と性格に関係がある」という結論が出る程度であれば、「血液型と性格に関係があるとは今のところ言えない」というのが、科学的態度です。

では、検定の結論は結局何を言っているのでしょうか？ それは、今日の問題の例で言えば、

私は、おじさんが『**50%**の確率で当たる』と言うのは間違いだ、と判断する。  
ただし、私は**100**回中**5**回はウソを言う(第1種の誤りを犯す)人間である。  
今回私がウソを言っているかどうか、それは誰にもわからない。

と言っているのに等しいのです。

この程度のことしか言っていないのに、検定にはどういう意味があるのでしょうか？

それは、検定とは、少ない数のデータしか、しかも1度だけしか調べられないときに、「それだけのデータからでも十分な確信をもって述べられる疑いだけを述べる」方法ということなのです。何度も検定できるほどデータがあつまるのなら、検定を用いるのは不適切です。

くじびぎの場合は、今くじをひいた数だけのデータしか集められませんから、検定を用いるのは妥当です。しかし、上で述べた血液型と性格の問題では、長い間に集められたたくさんのデータについて検定を用いるのは、検定の目的からはずれています。この問題については、第12回の講義で説明します。