

III-2. 確率分布

III-2-1. 確率の計算と二項分布モデル

頻度主義の立場に立つと、はじめに、正確に確率どおりに現れる理想的な現象がどんな形（確率分布）を描くのかを考えなければなりません。そのために、例として出されるのが、理想的なコインとか理想的なサイコロです（コインやサイコロなど賭け事が理想的だと言ったものではありません。念のために）。 n 回コインを投げて k 回表になる。あるいは、 n 回サイコロを投げて k 回1の目が出る。その確率を $P(k)$ としたとき（ k となる確率 Probability という感じでこの記号にしました。）、 0 から n までの k の値について $P(k)$ を計算することができます。 k を横軸にして $P(k)$ の値をとったのが確率分布です。 n 回コインを投げて、何回表が出るか、もっともありそうな回数は何回かと訊かれたら、たいていの人が $\frac{n}{2}$ 回と答えます。これを k の期待値と言います。この場合は、直観的に、平均値にも中央値にもなっていることがわかります。 k の分布は 0 から n まで広がっていますから、分布の中心でここが一番確率が高いというわけです。サイコロで1が出る回数の場合は期待値は $\frac{n}{6}$ で、分布が左右相称ではないので、真ん中という感じはしないかもしれませんが、ここがもっとも確率が高く、平均値かつ中央値になっていて、分布の中心です。このような分布を二項分布と言います。表が出ることの排反事象（表でないこと）は裏が出ることで、1の目が出ることの排反事象（1でないこと）が1以外の目が出ることで二項対立的だからです。

実際に n 回コインを投げるかサイコロを振るかして、 x 回、表ないし1の目がでたら、その値を k の値を表す横軸で見つけて、その点での確率 $P(x)$ が、理想的な確率分布上で x となる確率です。何故、二項分布で考えるのかというと、それ以外にうまいモデルが思い浮かばないからです。ただそれだけのことです。

次に、回数という不連続な値の確率に関するモデルを拡張して、背の高さとか、光の強さとか、値段とか、連続的な値にこの考え方に使える確率モデルを理論的につくります。さらに、そのような確率モデルで説明できる複数の確率変数の差や比なども確率変数と考えられるので、その変数の確率的変動を説明する確率モデルも作ります。それらのモデルの実際の形（平均値・ばらつき方）の推定のしかたを考えて、2章で説明した頻度主義の考え方で実際の判断（検定）をします。

この流れに従って、まず、二項分布モデルを数式として作るための作業をします。「組み合わせ」の数とか、積集合、和集合、確率の演算を説明するために、集合の概念とその演算の知識が必要になるので、使われる演算子の記述の仕方を含めて、基礎的な説明をまとめてします。

集合とはいくつかの要素をまとめたものです。普通は $\{ \}$ で表します。たとえば、 $1,3,4,6,7$ を要素としてもつ数のまとまりを一つの集合として表せば、 $\{1,3,4,6,7\}$ です。 $\{\text{犬}\}$ ならば、すべての犬を表します。隣の犬も、自分のうちの犬も、セントバーナードも、チワワも要素として含まれます。 $\{\text{動物}\}$ ならば、犬や猫も要素として含まれていて、それらもまた集合として表せます。また、その集合の要素は、すべて $\{\text{動物}\}$ の要素の中に含まれます。そのよ

うなものを部分集合と言います。ある集合の中に部分集合があれば、 $\{A\} \subset \{B\}$ のように、 \subset で表します。 $\{A\}$ は $\{B\}$ に部分集合として含まれるということです。 $\{犬\} \subset \{動物\}$ という風に表せます。これを包含関係と言います。 $\{A\} \subset \{B\} \subset \{C\}$ ならば、 $\{A\} \subset \{C\}$ であることは直感的に明らかです。ある要素がある集合に含まれるとき、 $e \in A$ と \in を使って表します。「 e は集合 $\{A\}$ の要素だ。」あるいは「集合 $\{A\}$ は要素として e を含む。」ということです。 $\{隣の犬\} \in \{犬\}$ です。ある集合 $\{A\}$ に含まれるがそれとは違う集め方をした別の集合 $\{B\}$ にも含まれるという場合、 $A \cap B$ と \cap で表します。積集合と言い、 A かつ B と読みます。たとえば、 $A = \{犬\}$ で $B = \{雌\}$ ならば $A \cap B = \{雌犬\}$ です。 A である要素のすべてと、 B である要素の全てを集めたものを和集合 $A \cup B$ と言います。 $A = \{犬\}$ で $B = \{雌\}$ ならば $A \cup B = \{人間を含めたすべての動物のメスと犬のすべて\}$ です。ですから、この中に部分集合として $\{雌犬\}$ が含まれます。 $(A \cup B) \supset (A \cap B)$ です。(なんだか、ウィーメエンズリブの人に怒られそうな文章になってしまいました。悪意があったわけではありません。犬が好きで、女性が好きなので、たまたま興味があるものを2つあげたらこうなっていました。)

記号であらわすと、

$$A = \{a_1, a_2 \dots a_l, c_1, c_2, \dots c_n\} \text{ and } B = \{b_1, b_2, \dots b_m, c_1, c_2, \dots c_n\},$$

の場合

$$D = A \cup B = \{a_1, a_2 \dots a_l, b_1, b_2, \dots b_m, c_1, c_2, \dots c_m\}$$

式 1

$$C = A \cap B = \{c_1, c_2, \dots c_m\}$$

式 2

A が起きて、その条件のもとに B が起きることを $B|A$ と表します。 $(A \cap B) = (B \cap A)$ ですが、 $(A|B) \neq (B|A)$ です。これらを使うと、何回かサイコロを振って出てくるもの組み合わせを集合としてとらえて、その組み合わせの数やそのような組み合わせになる確率を考えることができます。たとえばサイコロを2回振って目が1とそれ以外の数になる組み合わせを考えます。 A が1になる。 B が1以外になることにします。組み合わせの数を N としてその組み合わせになる確率を P とします。1回の試行であることが起きる確率を p とします。それが起こらない確率は q で $p + q = 1$ です。コインの場合は表になる確率 $p = \frac{1}{2}$ 、表にならない確率 $q = \frac{1}{2}$ 、サイコロの場合は1の目が出る確率 $p = \frac{1}{6}$ 、1の目が出ない確率 $q = \frac{5}{6}$ ということです。互いに排反事象ですから $p + q = 1$ ということは理解できますね。

1回サイコロを振る場合は、 A か B しかないので

$$N(A) = 1, N(B) = 1, N(A \cup B) = 2, P(A) = \frac{1}{6}, P(B) = \frac{5}{6}, P(A \cup B) = 1,$$

ですね。

サイコロを2回振る場合は、 $A|A, B|A, A|B, B|B$ の四通りがあって、これらは互いに排反事象(あることが起きた場合には他のことは起こらない)です。

A と B が背反事象ならば

$$P(A \cup B) = P(A) + P(B)$$

式 3

$$N(A|A) = 1, \quad N(B|A) = 1, N(A|B) = 1, N(B|B) = 1, N(A|A) \cup (B|A) \cup (A|B) \cup (B|B) = 4$$

それぞれの確率は

$$P(A|A) = \frac{1}{6} \cdot \frac{1}{6}, \quad P(B|A) = \frac{1}{6} \cdot \frac{5}{6}, P(A|B) = \frac{5}{6} \cdot \frac{1}{6}, P(B|B) = \frac{5}{6} \cdot \frac{5}{6}$$

$$P((A|A) \cup (B|A) \cup (A|B) \cup (B|B)) = 1$$

B|AはAが起きたという条件のもとにBが起きるのですから、与えられた条件を式にすると

$$P(B|A) = P(A)P(B|A)$$

と書くべきですが、この場合は、AにかかわりなくBがおきるので、P(B)は一定で

$$P(B|A) = P(B)$$

です。たがいに独立したあること(A)とあること(B)が同時に起きる確率はそれら個々の確率の積ですね。この関係を式で表すと次のようになります。

$$P(A \cap B) = P(A)P(B)$$

式 4

したがって

$$P(B|A) = P(A)P(B|A) = P(A)P(B)$$

式 5

と書けます。

ところで、 $P(A \cap B)$ はAとBが同時に起きるという意味ですが、実際の時間の中で「同時」ということではなくて、数学的に考えた場合に「同時」ととらえられる、一つの試行の中でということです。「同時」ということよりは「互いに独立した」という条件の方が重要なのです。

組み合わせの数を整理すると

$$AA = \{A \text{ が 2 回}\}$$

$$AB = \{A \text{ が 1 回}\}$$

$$BB = \{A \text{ が 0 回}\}$$

と整理できて、

$$AA = A|A$$

$$AB = (B|A) \cup (A|B)$$

$$BB = B|B$$

$$N(AA) = 1, \quad N(AB) = 2, N(B|B) = 1, N(AA) \cup (AB) \cup (BB) = 4$$

1 試行について 2 通り、そのそれぞれについて次の試行で 2 通りの結果があるのだから、

その2つの試行を1セットとして、1試行と考えれば4通りあるということになります。

1試について n_1 とおりの結果があり次の試行でそのそれぞれに結果について n_2 通りの結果があれば、2回の試行を1セットの試行と考えれば、あらわれる結果には $n_1 \cdot n_2$ とおりの組み合わせがあります。つまり、2つを組み合わせた場合の数は一つひとつの場合の数の掛け算なのです。

実際に、考えてみます。サイコロを振って、1の目が出るのがA、それ以外の目が出るのをBとします。4回投げると1試行だとすると、一回目にサイコロの目が1で、残りの3回が1以外になるのは、(|B|B|B|A)と表されて、

$$\begin{aligned} P(|B|B|B|A) &= P(A \cap B \cap B \cap B) = P(A)P(B)P(B)P(B) \\ &= pqqq = \frac{1555}{6666} = \frac{125}{1296} \end{aligned}$$

ここで、サイコロの目が1回だけ1で、後の3回は1以外という組み合わせをすべて書くと、(|B|B|B|A)、(|B|B|A|B)、(|B|A|B|B)、(|A|B|B|B)の4通りがあります。

サイコロの目が1回だけ1で、後の3回は1以外になる確率は、 p^1q^3 の4倍、

$4pq^3 = 4\left(\frac{1}{6}\right)\left(\frac{5}{6}\right)^3 = \frac{500}{1296}$ になります。ここで $(p+q)^n$ という式の展開を考えます。

$$n=1 \text{ ならば、 } (p+q)^1 = p+q$$

$$n=2 \text{ ならば } (p+q)^2 = p^2 + 2pq + q^2$$

$$n=3 \text{ ならば } (p+q)^3 = p^3 + 3p^2q + 3pq^2 + q^3$$

$$n=4 \text{ ならば } (p+q)^4 = p^4 + 4p^3q + 6p^2q^2 + 4pq^3 + q^4$$

n=4の場合の式の後ろから2項目を見てください。サイコロの目が1回だけ1で、後の3回は1以外になる確率は、 p^1q^3 の4倍、 $4pq^3$ と同じになっています。式を展開していく過程で生じる $pqqq$ 、 $qpqq$ 、 $qqpq$ 、 $qqqp$ の4つの項を pq^3 の項として一つにまとめたのですから、当たり前だといわれれば、それまでです。しかし、そうだとすれば、n回サイコロを振って、1の目がk回出る確率は、後ろからk+1番目の項だということになります。前から数えると、n-(k-1)番目の項です。一番先頭が、n回振って、n回1の目が出る(k=nになる)確率で、一番最後が一回も1の目が出ない(k=0になる)確率です。p+q=1だから、それを何乗しても1で、常に確率の総和は1だということも確認できます。各項のkの値を横軸にして、それぞれの確率を縦軸にプロットしたものを二項分布といいます。各項の前の部分に数字がありますが、この数字を二項係数と言いい次のような記号で表します。

$$\binom{n}{k}$$

$(p+q)^n$ を展開したときの前からn-(k-1)番目の項の係数という意味です。

$$(p + q)^4 = \binom{4}{4}p^4 + \binom{4}{3}p^3q + \binom{4}{2}p^2q^2 + \binom{4}{1}pq^3 + \binom{4}{0}q^4$$

と表せます。一般化して書くと

$$(p + q)^n = \sum_{k=0}^n \binom{n}{k} p^{n-k} q^k$$

となります、

k なる確率についてだけ取り出して書くと

$$P(k) = \binom{n}{k} p^{n-k} q^k$$

ですが、 $q = 1 - p$ なので、 p だけの式に書き換えると

$$P(k) = \binom{n}{k} p^{n-k} (1 - p)^k$$

となります。1回の確率が p であること n 回繰り返した時の二項分布を

$$B(n, p)$$

と書きます。

さて、残された問題は $\binom{n}{k}$ をどのように計算するかです。式1は

$$(p + q)^n = \sum_{k=0}^n \binom{n}{k} p^{n-k} q^k$$

となっていて p と q のべき乗数の和 $(n - k) + k$ はどの項も n です。 $ppppqqq$ のように書くと、 n この文字が列を作っています。もしこれが a, b, c, d, e, f のように n 個の文字が並んでいて、その順番を入れ替えたときに、何個の並び方があるかを考えます。最初に来る可能性があるのは n 個あって、その次は先頭になったもの以外、その次は先頭とその次目以外の文字が並ぶと考えていけば、それぞれについて、組み合わせの数は、 $n, n - 1, n - 2 \dots$ のようになって、最後は一つだけになります。独立した組み合わせ同士を組み合わせた組み合わせの数は、組み合わせ数同士の掛け算ですから、組み合わせの総数は $n(n - 1)(n - 2) \dots 3 \cdot 2 \cdot 1$ つまり、 $n!$ になります。さて、文字が n 個あるときに、 k 個を p 組、 $n - k$ 個を q 組に分けたとします。同じように考えると、 p 組だけの並び方の総数は $k!$ で q 組だけの並び方の総数は $(n - k)!$ です。組み合わせ数同士の掛け算が、組み合わせた場合の組み合わせ数の総数なので、

p 組の内部の組み合わせ数 \times q 組の内部の組み合わせ数 \times 組内の組み合わせを考慮しない p と q の組の違いだけの組み合わせ数 = 総組み合わせ数

$$k! (n - k)! \binom{n}{k} = n!$$

となります。これを変形すると、

$$\begin{array}{ccccccc}
& & & & 1 & & & & \\
& & & & 1 & & 1 & & \\
& & & 1 & & 2 & & 1 & \\
& & 1 & & 3 & & 3 & & 1 \\
& 1 & & 4 & & 6 & & 4 & & 1 \\
1 & & 5 & & 10 & & 10 & & 5 & & 1 \\
1 & & 6 & & 15 & & 20 & & 15 & & 6 & & 1
\end{array}$$

全体として以下の式の様に対称形になりますが

$$\binom{n}{k} = \binom{n}{n-k}$$

上下にみると次のような関係があります。

$$\binom{n}{k} = \binom{n-k}{k-1} + \binom{n-k}{k}$$

$$\begin{array}{ccccccc}
& & & & 1 & & & & \\
& & & & 1 & & 1 & & \\
& & & 1 & & 2 & & 1 & \\
& & 1 & & 3 & & 3 & & 1 \\
1 & & 4 & & 6 & & 4 & & 1 \\
1 & & 5 & & 10 & & 10 & & 5 & & 1 \\
1 & & 6 & & 15 & & 20 & & 15 & & 6 & & 1
\end{array}$$

$$2 = 1 + 1$$

$$3 = 1 + 2$$

$$4 = 1 + 3, \quad 6 = 3 + 3$$

1回の試行である事象が起こる確率を p として、それらを n 回繰り返すことを、二項分布の記号として、 $B(n, p)$ と表します。この時、ある事象が k 回起きる確率を $p(k)$ とすると、 $p(k)$ は二項分布にしたがうといくことで、次のように表します。

$$p(k) \sim B(n, p)$$

図5には、サイコロを振って何回1が出るかという確率 $p(k) \sim B\left(n, \frac{1}{6}\right)$ を ($1 \leq n \leq 10$)

について、図示しました。 n の増加によって、二項分布が、次第にシンメトリックな、正規分布に近づいていくことがおぼろげにわかります。

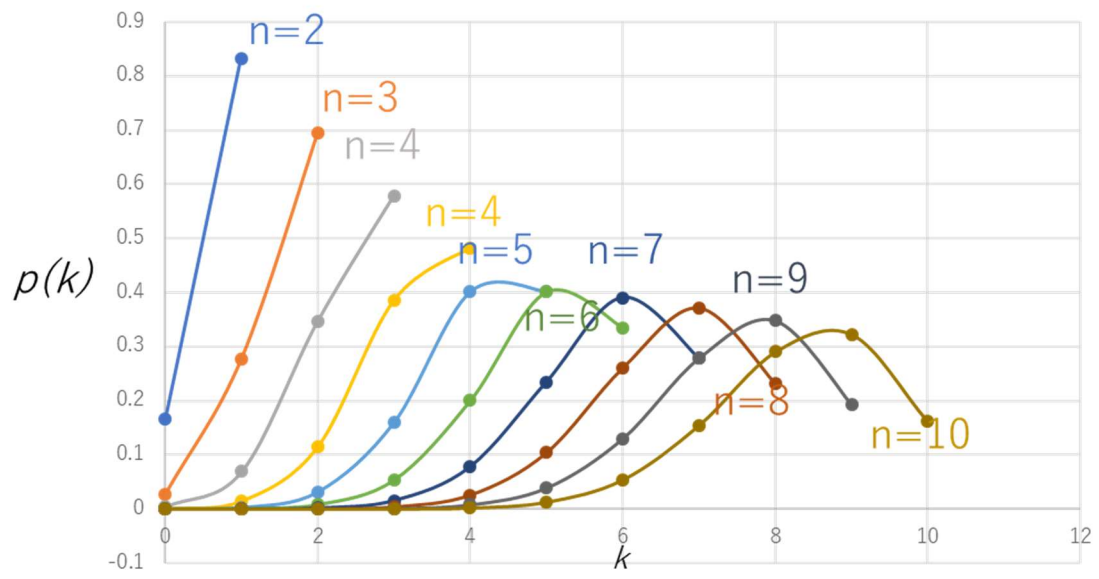


图. 5. 二項分布