

## ベイズ推定（基本編）

### 1 ベイズの定理

確率に関するベイズの定理を復習する。

#### 定義（条件付き確率）

事象 $A_i$ と $B_j$ を考え、 $A_i$ と $B_j$ が同時に起こる確率を $p(A_i, B_j)$ とする。そして、

$$p(B_j|A_i) = \frac{p(A_i, B_j)}{p(A_i)}$$

とする。 $p(B_j|A_i)$ を、 $A_i$ という条件のもとで $B_j$ が起こる条件付き確率をという。

この定義より、

$$p(A_i, B_j) = p(B_j|A_i)p(A_i)$$

がいえ、同様に、

$$p(A_i, B_j) = p(A_i|B_j)p(B_j)$$

もいえる。したがって、次の定理が得られる。

$$p(A_i|B_j) = \frac{p(B_j|A_i)p(A_i)}{p(B_j)} \quad (1)$$

を得る。このとき、 $P(A_i)$ を $A_i$ の事前確率、 $P(A_i|B_j)$ を $A_i$ の事後確率という。また、 $p(B_j|A_i)$ を尤度<sup>ゆうど</sup>という。

さて、

$$\sum_{i=1}^n \sum_{j=1}^m p(A_i, B_j) = 1$$

とする。このとき、

$$p(B_j) = \sum_{i=1}^n p(A_i, B_j) = \sum_{i=1}^n p(B_j|A_i)p(A_i) \quad (2)$$

であり、(1)と(2)より、以下の定理が得られる。

#### 確率に関するベイズの定理

$$\sum_{i=1}^n \sum_{j=1}^m p(A_i, B_j) = 1 \text{ ならば } p(A_i|B_j) = \frac{p(B_j|A_i)p(A_i)}{\sum_{i=1}^n p(B_j|A_i)p(A_i)} \quad (3)$$

**【例題 1】** ある工場の製品は、機械 $M_1$ では単位時間当たり 15 個、機械 $M_2$ では単位時間当たり 20 個の速さで製造されている。しかし、 $M_1$ では3%の不良品が、 $M_2$ では5%の不良品が含まれる。いま、一定の時間に $M_1$ と $M_2$ により製造された製品の中から任意に取り出した 1 個が不良品であった。この不良品が $M_1$ の製造品である確率を求めよ。

(解答) 単位時間において、事象 $A_1$ を $M_1$ の製造品、事象 $A_2$ を $M_2$ の製造品とする。さらに、事象 $B$ を不良品とする。このとき、

$$P(A_1) = \frac{15}{35}, \quad P(A_2) = \frac{20}{35}, \quad P(B|A_1) = \frac{3}{100}, \quad P(B|A_2) = \frac{5}{100}$$

である。これより、

$$p(B|A_1)p(A_1) = \frac{3}{100} \times \frac{15}{35} = \frac{9}{700}, \quad p(B|A_1)p(A_1) + p(B|A_2)p(A_2) = \frac{3}{100} \times \frac{15}{35} + \frac{5}{100} \times \frac{20}{35} = \frac{29}{700}$$

よって、求める確率である $P(A_1|B)$ は、ベイズの定理より

$$P(A_1|B) = \frac{p(B|A_1)p(A_1)}{p(B|A_1)p(A_1) + p(B|A_2)p(A_2)} = \frac{9}{29}$$

となる。□

## 2 ベイズ推定

事象 $B_j$ と $C_k$ は互いに独立、すなわち、 $p(B_j, C_k) = p(B_j)p(C_k)$ の場合には、

$$\begin{aligned} p(A_i, B_j, C_k) &= p(A_i|B_j, C_k)p(B_j, C_k) \\ p(A_i, B_j, C_k) &= p(B_j, C_k|A_i)p(A_i) \end{aligned}$$

が成り立つので、

$$p(A_i|B_j, C_k) = \frac{p(B_j, C_k|A_i)p(A_i)}{p(B_j, C_k)} = \frac{p(C_k|A_i)}{p(C_k)} \cdot \frac{p(B_j|A_i)p(A_i)}{p(B_j)} = \frac{p(C_k|A_i)p(A_i|B_j)}{p(C_k)} \quad (4)$$

が得られる。ここで、 $B_j$ と $C_k$ は互いに独立より、

$$p(C_k|A_i) = p(C_k|A_i, B_j) = \frac{p(A_i, B_j, C_k)}{p(A_i, B_j)} = \frac{p(A_i, C_k)p(B_j)}{p(A_i|B_j)p(B_j)} = \frac{p(A_i, C_k)}{p(A_i|B_j)}$$

であることに注意すると、

$$p(C_k) = \sum_{i=1}^n p(A_i, C_k) = \sum_{i=1}^n p(C_k|A_i)p(A_i|B_j)$$

となる。よって、(4)より次の定理を得る。

**定理 (ベイズ推定, またはベイズ更新)** 事象 $B_j$ と $C_k$ は互いに独立ならば、

$$p(A_i|B_j, C_k) = \frac{p(C_k|A_i)p(A_i|B_j)}{\sum_{i=1}^n p(C_k|A_i)p(A_i|B_j)} \quad (5)$$

[例題 4.1] 迷惑メールフィルターは、迷惑メール $A_1$ か一般メール $A_2$ を振り分けるものである。メール $A$ には「無料」、「未納料金」、「サイト利用料金」という言葉が含まれている。迷惑メールと一般メールに、それらが含まれる確率は、これまでのデータベースから、以下と表のとおりである。また、データベースから現時点での迷惑メールの割合は6割であることもわかっている。

	迷惑メール $A_1$	一般メール $A_2$
無料	0.12	0.02
未納料金	0.15	0.01
サイト利用料金	0.25	0.02

無料という言葉が含まれている事象を $B_1$ 、そうでない事象を $B_2$ とし、未納料金という言葉が含まれている事象を $C_1$ 、そうでない事象を $C_2$ とし、サイト利用料金という言葉が含まれている事象を $D_1$ 、そうでない事象を $D_2$ とする。このとき、 $p(A_1|B_1, C_1, D_1)$ をベイズ更新から推定せよ。

(解答) 表より、 $p(B_1|A_1) = 0.12$ である。また、 $p(A_1) = 0.6$ である。(4.3)より、

$$p(A_1|B_1) = \frac{p(B_1|A_1)p(A_1)}{p(B_1|A_1)p(A_1) + p(B_1|A_2)p(A_2)} = \frac{0.12 \times 0.6}{0.12 \times 0.6 + 0.02 \times 0.4} = 0.9000$$

である。つぎに、 $p(C_1|A_1) = 0.15$ であることから、(4.5)より

$$p(A_1|B_1, C_1) = \frac{p(C_1|A_1)p(A_1|B_1)}{p(C_1|A_1)p(A_1|B_1) + p(C_1|A_2)p(A_2|B_1)} = \frac{0.15 \times 0.9}{0.15 \times 0.9 + 0.01 \times 0.1} = 0.9926$$

となる。最後に、 $p(D_1|A_1) = 0.25$ であることから、(4.5)より

$$p(A_1|B_1, C_1, D_1) = \frac{p(D_1|A_1)p(A_1|B_1, C_1)}{p(D_1|A_1)p(A_1|B_1, C_1) + p(D_1|A_2)p(A_2|B_1, C_1)} = \frac{0.25 \times 0.9926}{0.25 \times 0.9926 + 0.02 \times 0.0074} = 0.9993$$

以上より、メール $A$ は迷惑メールである可能性が極めて高い。□

### 3 頻度確率とベイズ確率の考え方の違い

現代の統計学では、大きく頻度確率とベイズ確率と呼ばれる2種類の考え方で扱われる。以下、簡単にそれらを説明する。

統計を頻度確率で扱おうとする考え方は、できるだけ偏りのない観測を多数得ることで、公平・公正な判断を導くことに重きを置くというものである。したがって、多少手間をかけても良いから、より精確なデータ解析を必要とする用途に適している。しかし、欲しい情報を欲しいだけ観測できるかどうかは問題である。

統計をベイズ確率で扱おうとする考え方は、多少データが偏っていてもよいから、少ない観測で効率的に推定が行えることを優先するものである。そのためなら、外部の情報を事前知識として持ち込むことを厭わず、使えるものは使いたいとする考え方である。しかし、事前知識が不正確だと推定は怪しくなるというリスクは必然である。それでもそのリスクを負うことを承知してでも、推定値を効率的に得ようとする考え方である。具体的には、観測自体ができない状態を扱いたい場合、ある種の仮説を立てて、演繹的な推論を行うのである。

#### 4 確率分布に関するベイズの定理とベイズ推定

$\theta$ と $x$ を確率変数とする.  $f(\theta, x)$ を $\theta$ と $x$ に関する確率密度関数とする. さらに,  $f(x|\theta)$ を $\theta$ を固定したときの $x$ に関する確率密度関数とする. 確率に関する

$$p(A_i, B_j) = p(B_j|A_i)p(A_i), \quad p(A_i, B_j) = p(A_i|B_j)p(B_j)$$

から,  $f(\theta, x)$ と $f(x|\theta)$ との間に

$$f(\theta, x) = f(x|\theta)f(\theta), \quad f(\theta, x) = f(\theta|x)f(x) \tag{6}$$

という関係式を定義する. これより,

$$f(\theta|x) = \frac{f(x|\theta)f(\theta)}{f(x)} \tag{7}$$

を得る.  $f(\theta)$ を $\theta$ の事前確率分布,  $f(\theta|x)$ を $\theta$ の事後確率分布という. また,  $f(x|\theta)$ を $x$ の尤度関数という.

また, 確率に関する式

$$p(B_j) = \sum_{i=1}^n p(B_j|A_i)p(A_i)$$

に対応させて

$$f(x) = \int_{-\infty}^{\infty} f(x|\theta)f(\theta)d\theta \tag{8}$$

を定義する.

(8)と(9)から, 確率分布に関するベイズの定理

$$f(\theta|x) = \frac{f(x|\theta)f(\theta)}{\int_{-\infty}^{\infty} f(x|\theta)f(\theta)d\theta} \tag{9}$$

が得られる.  $x_1, x_2$ が互いに独立であるとき, ベイズ更新

$$p(A_i|B_j, C_k) = \frac{p(C_k|A_i)p(A_i|B_j)}{\sum_{i=1}^n p(C_k|A_i)p(A_i|B_j)}$$

に対応して,

$$f(\theta|x_1, x_2) = \frac{f(x_2|\theta)f(\theta|x_1)}{\int_{-\infty}^{\infty} f(x_2|\theta)f(\theta)d\theta} \tag{10}$$

となる. さらに,

$$k_2 = \frac{1}{\int_{-\infty}^{\infty} f(x_2|\theta)f(\theta)d\theta}$$

と置くと, (10)は

$$f(\theta|x_1, x_2) = k_2 f(x_2|\theta)f(\theta|x_1) \tag{11}$$

となる. これを一般化すると以下となる.

$$f(\theta|x^{(n)}) = k_n f(x_n|\theta)f(\theta|x^{(n-1)}), \quad \frac{1}{k_n} = f(x^{(n)}) = \int_{-\infty}^{\infty} f(x_n|\theta)f(\theta)d\theta \tag{12}$$

以上をまとめると, 以下となる.

確率分布に関するベイズ推定は

$$(1) \quad f(\theta|x_1) = k_1 f(x_1|\theta) f(\theta)$$

$$(2) \quad f(\theta|x_1, x_2) = k_2 f(x_2|\theta) f(\theta|x_1)$$

$$(3) \quad f(\theta|x_1, x_2, x_3) = k_3 f(x_3|\theta) f(\theta|x_1, x_2)$$

.....

$$(4) \quad f(\theta|x^{(n)}) = k_n f(x_n|\theta) f(\theta|x^{(n-1)}) \quad \text{ただし, } x^{(n)} = x_1, x_2, \dots, x_n$$

となる.

[例題 3] 表の出る確率が $\theta$ であるコインがある. このコインを3回投げたとき, 表, 表, 裏の順に出た. このときに $\theta$ に関する確率分布 $f(\theta|x^{(3)})$ を推定せよ. ただし, 事前分布は $f(\theta) = 1$ とせよ.

(解答) 1回目は表であるので $f(x_1|\theta) = \theta$ であり,  $f(\theta) = 1$ より

$$f(\theta|x_1) = k_1 f(x_1|\theta) f(\theta) = k_1 \theta$$

である.  $\int_0^1 f(\theta|x_1) d\theta = 1$ と

$$\int_0^1 k_1 \theta d\theta = \frac{k_1}{2} [\theta^2]_0^1 = \frac{k_1}{2}$$

より,  $k_1 = 2$ を得る. よって,

$$f(\theta|x_1) = 2\theta$$

となる. 次に, 2回目も表より $f(x_2|\theta) = \theta$ であり,  $f(\theta|x_1) = 2\theta$ より

$$f(\theta|x^{(2)}) = k_2 f(x_2|\theta) f(\theta|x_1) = 2k_2 \theta^2$$

であり,  $\int_0^1 f(\theta|x_2) d\theta = 1$ と

$$\int_0^1 2k_2 \theta^2 d\theta = \frac{2k_2}{3} [\theta^3]_0^1 = \frac{2k_2}{3}$$

より,  $k_2 = 3/2$ を得る. よって,

$$f(\theta|x^{(2)}) = \frac{3}{2} \cdot 2\theta \cdot \theta = 3\theta^2$$

となる. 最後に, 3回目は裏より $f(x_3|\theta) = 1 - \theta$ であり,  $f(\theta|x^{(2)}) = 3\theta^2$ より

$$f(\theta|x^{(3)}) = k_3 f(x_3|\theta) f(\theta|x^{(2)}) = 3k_3 \theta^2 (1 - \theta)$$

であり,  $\int_0^1 f(\theta|x_2) d\theta = 1$ と

$$\int_0^1 3k_3 \theta^2 (1 - \theta) d\theta = 3k_3 \left[ \frac{\theta^3}{3} - \frac{\theta^4}{4} \right]_0^1 = \frac{k_3}{4}$$

より,  $k_3 = 4$ を得る. よって,  $\theta$ に関する確率分布 $f(\theta|x^{(3)})$ は,

$$f(\theta|x^{(3)}) = 4 \cdot 3\theta^2 \cdot (1 - \theta) = 12\theta^2(1 - \theta)$$

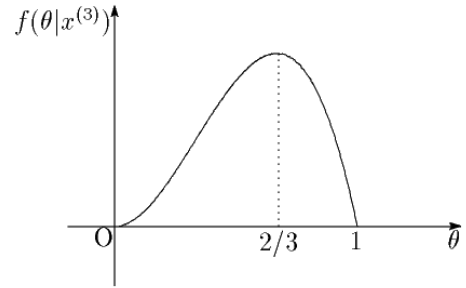
となる.

(解説)  $\theta$  関する確率分布  $f(\theta|x^{(3)}) = 12\theta^2(1-\theta)$  を考える.

$$f' = 12\theta(2-3\theta)$$

$0 < \theta < 1$  より,  $\theta = 2/3$  で  $f(\theta|x^{(3)})$  は極大値をとる. グラフは右図である.

そして, 3個の統計から得られたこのコインの表の出る確率は, 事前分布を一様分布として仮定した場合,  $\theta = 2/3$  が確率的に最も高いということを意味するのである.



### 練習問題

[1] ある国民の0.02%がある病気に罹患している. この病気のある検査方法では, 実際に病気に罹患している人が陽性と判定される確率は90%で, 病気に罹患していない人が陰性と判定される確率は80%である. ある人がこの病気の検査を受けて陽性であった. この人が病気に罹患している確率 (%) を, 小数点以下第3位まで求めよ.

[2] ある記憶デバイスは, A社, B社, C社で製造されている. このデバイスの市場のシェアは, それぞれ50%,30%,20%である. また, このデバイスが不良品として返品される確率は, それぞれ2%,5%,9%である. ある人がこのデバイスを購入したら不良品であった. このデバイスがA社の製品である確率 (%) を, 小数点以下第3位まで求めよ.

[3] 表の出る確率が $\theta$ であるコインがある. このコインを3回投げたとき, 裏, 表, 裏の順に出た. このときに $\theta$ 関する確率分布 $f(\theta|x^{(3)})$ を推定せよ. ただし, 事前分布は $f(\theta) = 1$ とせよ. また,  $f(\theta|x^{(3)})$ の極大値をとる $\theta$ を求めよ.