

統計的な推測【確率変数の分散と標準偏差】 p.58~60

1 学習内容の説明 ⇒ 2 問題演習 ⇒ 3 振り返り(確認テスト・相互採点・リフレクションの記入)

【態度目標】しゃべる、質問する、説明する、動く、協力する、貢献する

【内容目標】確率分布や期待値を用いて分散の値を求められるようになろう

□確率変数の分散と標準偏差

同じ期待値をもつ確率変数であっても、その分布は同じとは限らない。

値が期待値の近くに集中している分布もあれば、値が期待値から遠くに散らばっている分布もある。

ここでは、確率変数 X のとる値が、 X の期待値からどの程度散らばっているかを表す量について考えよう。

X の確率分布が右の表で与えられ、
その期待値が m であるとする。

$$\text{期待値} = \text{平均}$$

X	x_1	x_2	……	x_n	計
P	p_1	p_2	……	p_n	1

このとき、 X の各値と m とのへだたりの程度を表す量として

$$(x_1 - m)^2, (x_2 - m)^2, (x_3 - m)^2, \dots, (x_n - m)^2$$

が考えられ、 $(X - m)^2$ はこれらの値をとる確率変数である。

確率変数 $(X - m)^2$ の期待値 $E((X - m)^2)$ を、
確率変数 X の 分散 といい、 $V(X)$ で表す。

$\overbrace{V(X)}$ の V は、分散を意味する英語
variance の頭文字である。

分散は、「確率分布のばらつきを示すもの」

これは
(平均との差) の2乗
つまり
(偏差) の2乗
この平均を求めれば
分散を求めることができる
(基本数 Iと同じ！)

すなわち、 $V(X)$ は次の式で表される。

$$V(X) = (x_1 - m)^2 p_1 + (x_2 - m)^2 p_2 + \dots + (x_n - m)^2 p_n$$

また、和の記号 Σ を使って表すと、次のようになる。

$$V(X) = \sum_{k=1}^n (x_k - m)^2 p_k \quad \dots \quad ①$$

確率変数の分散

$$\begin{aligned} V(X) &= E((X - m)^2) \\ &= (x_1 - m)^2 p_1 + (x_2 - m)^2 p_2 + \dots + (x_n - m)^2 p_n \\ &= \sum_{k=1}^n (x_k - m)^2 p_k \end{aligned}$$

偏差の2乗の期待値 (平均)

例5) 2枚の硬貨を同時に投げるとき、表が出る硬貨の枚数 X について、

期待値 m は、 $m = 1$ である。

よって、 X の分散は

$$V(X) = (0 - 1)^2 \cdot \frac{1}{4} + (1 - 1)^2 \cdot \frac{2}{4} + (2 - 1)^2 \cdot \frac{1}{4} = \frac{1}{2} \quad \text{総}$$

X	0	1	2	
$(X - m)^2$	1	0	1	計
P	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$	1

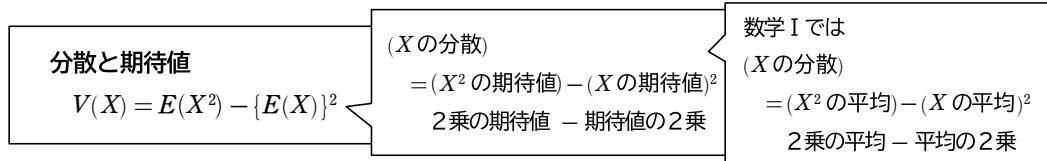
p_k が規則性がないので地道に計算

X の分散を表す前ページの式①の右辺を変形してみよう。

$$\begin{aligned}
 V(X) &= \sum_{k=1}^n (x_k - m)^2 p_k \\
 &= \sum_{k=1}^n (x_k^2 - 2mx_k + m^2) p_k \\
 &= \sum_{k=1}^n x_k^2 p_k - 2m \sum_{k=1}^n x_k p_k + m^2 \sum_{k=1}^n p_k \\
 &= E(X^2) - 2m \cdot m + m^2 \cdot 1 \\
 &= E(X^2) - m^2
 \end{aligned}$$

$$\begin{aligned}
 \sum_{k=1}^n x_k p_k &= m \\
 \sum_{k=1}^n p_k &= 1
 \end{aligned}$$

ここで、 $m = E(X)$ であるから、次の等式が成り立つ。



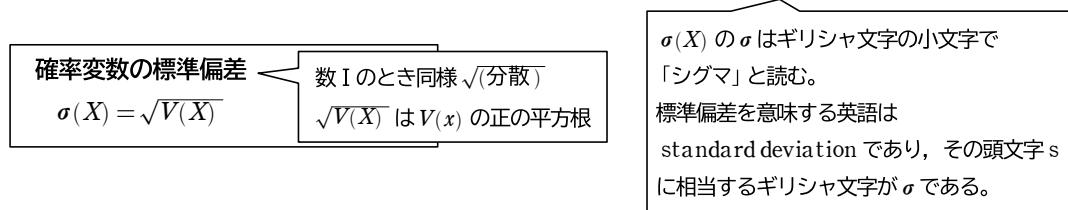
例6) 1個のさいころを投げて出る目 X の分散 $V(X)$ を求める。

$$56 \text{ ページの例2により } E(X) = \frac{7}{2}, \quad \boxed{\text{期待値}}$$

$$57 \text{ ページの例4により } E(X^2) = \frac{91}{6} \quad \boxed{\text{2乗の期待値}}$$

$$\text{であるから } V(X) = E(X^2) - [E(X)]^2 = \frac{91}{6} - \left(\frac{7}{2}\right)^2 = \frac{35}{12} \quad \text{総}$$

分散 $V(X)$ は確率変数 $(X - m)^2$ の期待値であるから、 X の測定単位が、たとえば cm であるとき、 $V(X)$ の単位は cm^2 となる。そこで、 X の測定単位と同じ単位である $\sqrt{V(X)}$ を散らばりの度合いを表す数値として用いることが多い。 $\sqrt{V(X)}$ を X の **標準偏差** といい、 $\sigma(X)$ で表す。



例7) 1個のさいころを投げて出る目 X の標準偏差 $\sigma(X)$ を求める。

$$\text{例6 により } V(X) = \frac{35}{12} \text{ であるから}$$

$$\sigma(X) = \sqrt{V(X)} = \sqrt{\frac{35}{12}} = \frac{\sqrt{105}}{6} \quad \text{総}$$

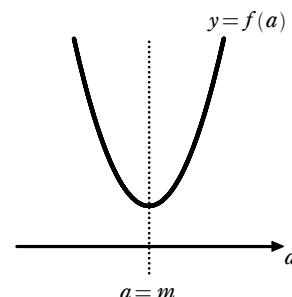
確率変数 X の期待値、分散、標準偏差を、それぞれ X の分布の **平均**、**分散**、**標準偏差** ともいう。標準偏差 $\sigma(X)$ は、 X の分布の平均 m を中心として、 X のとる値の散らばる傾向の程度を表している。標準偏差 $\sigma(X)$ の値が小さいほど、 X のとる値は、平均 m の近くに集中する傾向にある。

補足 分布の代表値の一つを a とし、分布のばらつきの関数を $f(a)$ を

$$f(a) = \frac{1}{n} \sum_{k=1}^n (x_k - a)^2 \text{ と定義する。ただし, } m \text{ は平均値とする。}$$

$$\begin{aligned} f(a) &= \frac{1}{n} \{ (x_1 - a)^2 + (x_2 - a)^2 + (x_3 - a)^2 + \dots + (x_n - a)^2 \} \\ &= \frac{1}{n} \{ na^2 - 2a(x_1 + x_2 + x_3 + \dots + x_n) + (x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2) \} \\ &= \frac{1}{n} \{ na^2 - 2anm + (x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2) \} \\ &= a^2 - 2am + \frac{1}{n} \cdot (x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2) \\ &= (a - m)^2 + \frac{1}{n} \cdot (x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2) - m^2 \end{aligned}$$

平方完成



よって $y = f(a)$ は、 $a = m$ で最小となり、
その値は「2乗の平均 引く 平均の2乗」になっている。

したがって、代表値の中でも平均 m を利用した方が、
ばらつき $f(a)$ を最小で評価できる。

【まとめ】

分布のばらつきは、代表値の中でも『平均を基準』に置いた方が最小で評価できる。
このとき、確率変数の平均と分散の関係が導ける。

$$V(X) = E(X^2) - [E(X)]^2$$

補足 分布の代表値の一つを a とし、分布のばらつきの関数を $f(a)$ を

$$f(a) = \frac{1}{n} \sum_{k=1}^n |x_k - a| \text{ と定義する。ただし, } m \text{ は平均値とする。}$$

$$f(a) = \frac{1}{n} \{ |x_1 - a| + |x_2 - a| + |x_3 - a| + \dots + |x_n - a| \}$$

ここで $x_1 \leq x_2 \leq x_3 \leq \dots \leq x_n$ として考える

区間 $(-\infty, x_1), [x_2, x_3], [x_3, x_4], \dots, [x_n, \infty)$ において、 $f(a)$ は a の1次関数になる

区間 $(-\infty, x_1)$ では、

$$\begin{aligned} f(a) &= \frac{1}{n} \{ -(a - x_1) - (a - x_2) - (a - x_3) - \dots - (a - x_n) \} \\ &= \frac{1}{n} \{ -na + (x_1 + x_2 + x_3 + \dots + x_n) \} \end{aligned}$$

$$= -a + \frac{1}{n}(x_1 + x_2 + x_3 + \dots + x_n)$$

$$= -a + m$$

区間 $[x_1, x_2]$ では、

$$\begin{aligned} f(a) &= \frac{1}{n}\{(a-x_1)-(a-x_2)-(a-x_3)-\dots-(a-x_n)\} \\ &= \frac{1}{n}\{-(n-2)a+(-x_1+x_2+x_3+\dots+x_n)\} \\ &= \left(-1 + \frac{2}{n}\right)a + \frac{1}{n}(-x_1+x_2+x_3+\dots+x_n) \end{aligned}$$

区間 $[x_2, x_3]$ では、

$$\begin{aligned} f(a) &= \frac{1}{n}\{(a-x_1)+(a-x_2)-(a-x_3)-\dots-(a-x_n)\} \\ &= \frac{1}{n}\{-(n-4)a+(-x_1-x_2+x_3+\dots+x_n)\} \\ &= \left(-1 + \frac{4}{n}\right)a + \frac{1}{n}(-x_1-x_2+x_3+\dots+x_n) \end{aligned}$$

区間 $[x_n, \infty)$ では、

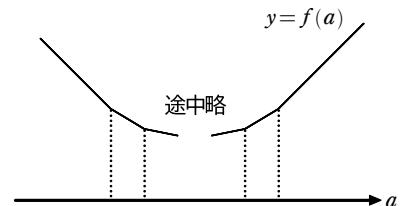
$$\begin{aligned} f(a) &= \frac{1}{n}\{(a-x_1)+(a-x_2)+(a-x_3)+\dots+(a-x_n)\} \\ &= \frac{1}{n}\{na-(x_1+x_2+x_3+\dots+x_n)\} \\ &= a - \frac{1}{n}(x_1+x_2+x_3+\dots+x_n) \\ &= a - m \end{aligned}$$

よって、ばらつき $y=f(a)$ は、各区間をつなぎ合わせた下に凸の折れ線となる。

$y=f(a)$ のグラフは、左側は傾き -1 の半直線で、右側は傾き $+1$ の半直線になっている

$f(a)$ の最小は、 n が奇数のとき、 $a=x_{\frac{n+1}{2}}$ のとき、 n が偶数のとき、 a が区間 $[x_{\frac{n}{2}}, x_{\frac{n}{2}+1}]$ にあるときであり

ばらつき $f(a)$ を最小にする a は、平均値ではなく「中央値」である。



【まとめ】

分散は、平均を基準にした方が、コンパクトでよい。

偏差の絶対値の平均は、データの中央値が最小となり、必ずしも平均値ではない。

ちなみに $\frac{1}{n} \sum_{k=1}^n |x_k - m|$ (m は平均) は、平均絶対偏差と呼ばれている。

(高校の統計学では、平均絶対偏差は扱いにくく、標準偏差が用いられている。)

また、ばらつき $f(a)$ を一般化して、 $f(a) = \frac{1}{n} \sum_{k=1}^n |x_k - a|^p$ と定義すれば、

平均絶対偏差は $p=1$ のとき、標準偏差は $p=2$ のときである。

この一般化された $f(a)$ に対して、 $\{nf(a)\}^{\frac{1}{p}}$ を L^p ノルムと言い、関数解析の分野で扱われている。