

第5章 データの分析

2 5 データの整理、データの代表値

① データ

運動の記録や所属クラスなどのように、ある特性を表す数量を **変量** といい、ある変量の測定値や観測値、調査結果などの集まりを **データ** という。

② 度数分布

**度数分布表** データのとり値をいくつかの区間に区切って階級を定め、各階級に度数を対応させた表。各階級の真ん中の値を **階級値** という。

**ヒストグラム** 度数分布表を柱状のグラフで表したもの。

③ 代表値

**平均値** 変量  $x$  のデータの値が  $x_1, x_2, \dots, x_n$  であるとき、このデータの平均値  $\bar{x}$  は 
$$\bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

**最頻値 (モード)** データにおいて、最も個数の多い値。データが度数分布表に整理されているときは、度数が最も大きい階級の階級値。

**中央値 (メジアン)** データを値の大きさの順に並べたとき、中央の位置にくる値。データの大きさが偶数のときは、中央の2つの値の平均値。

2 6 データの散らばりと四分位数

① 四分位数

**範囲** データの最大値から最小値を引いた差。

**四分位数** データの値を大きさの順に並べたとき、4等分する位置の値。小さい方から順に、**第1四分位数**、**第2四分位数**、**第3四分位数** といい、順に  $Q_1$ 、 $Q_2$ 、 $Q_3$  で表す。第2四分位数はデータの中央値である。

**四分位範囲** データの第3四分位数  $Q_3$  から第1四分位数  $Q_1$  を引いた差  $Q_3 - Q_1$

② 箱ひげ図

**箱ひげ図** データの分布を、次の5つの値で表した図。平均値を記入することもある。最小値、第1四分位数 ( $Q_1$ )、中央値 ( $Q_2$ )、第3四分位数 ( $Q_3$ )、最大値



**外れ値** データの中に、他の値から極端に離れた値が含まれるとき、そのような値を **外れ値** という。外れ値の基準として、たとえば次のようなものがある。  
(第1四分位数  $-1.5 \times$  四分位範囲) 以下の値  
(第3四分位数  $+1.5 \times$  四分位範囲) 以上の値

2 7 分散と標準偏差

① 分散、標準偏差

変量  $x$  のデータの値を  $x_1, x_2, \dots, x_n$ 、その平均値を  $\bar{x}$  とする。

**偏差** 変量  $x$  のデータの各値から平均値  $\bar{x}$  を引いた差  $x_1 - \bar{x}$ 、 $x_2 - \bar{x}$ 、 $\dots$ 、 $x_n - \bar{x}$  これを  $x - \bar{x}$  で表す。

**分散** 偏差の2乗の平均値。 $s^2$  で表す。

$$s^2 = \frac{1}{n} \{ (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \}$$

$$(x \text{ のデータの分散}) = (x^2 \text{ のデータの平均値}) - (x \text{ のデータの平均値})^2$$

**標準偏差** 分散の正の平方根。 $s$  で表す。

$$s = \sqrt{\text{分散}} = \sqrt{\frac{1}{n} \{ (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \}}$$

$$(x \text{ のデータの標準偏差}) = \sqrt{(x^2 \text{ のデータの平均値}) - (x \text{ のデータの平均値})^2}$$

研究 変量の変換

$a$ 、 $b$  は定数とする。変量  $x$  のデータから  $y = ax + b$  によって新しい変量  $y$  のデータが得られるとき、 $x$ 、 $y$  のデータの平均値を  $\bar{x}$ 、 $\bar{y}$ 、分散を  $s_x^2$ 、 $s_y^2$ 、標準偏差を  $s_x$ 、 $s_y$  とすると 
$$\bar{y} = a\bar{x} + b, \quad s_y^2 = a^2 s_x^2, \quad s_y = |a| s_x$$

2 8 2つの変量の間の関係

① 散布図

**散布図** 2つの変量からなるデータを点として平面上に図示したもの。

**正の相関** 2つの変量からなるデータにおいて、一方が増えると他方も増える傾向がみられるとき、2つの変量の間には正の相関があるという。

**負の相関** 2つの変量からなるデータにおいて、一方が増えると他方が減る傾向がみられるとき、2つの変量の間には負の相関があるという。

注 正の相関と負の相関のいずれの傾向もみられないときは、2つの変量の間には相関がないという。

② 相関係数

**共分散**  $x$  の偏差と  $y$  の偏差の積の平均値。 $s_{xy}$  で表す。

$$s_{xy} = \frac{1}{n} \{ (x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}) \}$$

**相関係数** 相関の正負と強弱を表す値。 $r$  で表す。

$$r = \frac{s_{xy}}{s_x s_y} = \frac{\frac{1}{n} \{ (x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}) \}}{\sqrt{\frac{1}{n} \{ (x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \}} \sqrt{\frac{1}{n} \{ (y_1 - \bar{y})^2 + \dots + (y_n - \bar{y})^2 \}}} = \frac{(x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})}{\sqrt{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2} \sqrt{(y_1 - \bar{y})^2 + \dots + (y_n - \bar{y})^2}}$$

相関係数  $r$  は  $-1 \leq r \leq 1$  であり、 $r$  が1に近いほど正の相関が強く、 $-1$ に近いほど負の相関が強い。相関がないとき、 $r$  は0に近い値をとる。

注 相関係数は、外れ値の影響を受けやすい値である。

③ 相関関係と因果関係

一般に、2つの変量の間に相関関係があるからといって、必ずしも因果関係があるとはいえない。

④ 質的データをとる2つの変量の間の関係

**量的データ** 身長や体重など、数値として得られるデータ。

**質的データ** 所属クラスや都道府県など、数値ではないものとして得られるデータ。

**分割表** たとえば、合否が判定されるある試験において、受験者100人全員を対象に、教材Aを使用して学習したか調べるとする。2つの変量(合否、使用したかどうか)の間の関係を調べて、右の図のようにまとめたとき、この表を **分割表** または **クロス集計表** という。

	合	否	計
Aの使用：有	9	5	14
Aの使用：無	42	44	86
計	51	49	100

2 9 仮説検定の考え方

① 仮説検定の考え方

得られたデータをもとに、ある主張[1]が正しいかどうかを判断する、次のような手法を **仮説検定** という。

- ① 主張[1]と反する仮定を立てる。(主張[2]とする。)
- ② 主張[2]のもとで、実際に起こった出来事が起こりにくい出来事かどうかを調べる。
- ③ ②で調べた結果、実際に起こった出来事は十分起こりにくいと判断するとき、主張[2]の仮定は正しくないと判断できる。
- ④ 主張[1]は正しいと判断してもよいと考えられる。

補足 ③で、実際に起こった出来事が十分起こりにくいと判断しないときは、主張[2]の仮定は否定できず、主張[1]は正しいと判断できない。このとき、主張[2]が正しいと判断できるわけではない。

発展 仮説検定の考え方と反復試行の確率

① 組合せ

$n$  個から  $r$  個取る組合せ(異なる  $n$  個のものから異なる  $r$  個を取り出して作る組合せ)の

総数は 
$${}_n C_r = \frac{n(n-1) \cdots (n-r+1)}{r(r-1) \cdots 3 \cdot 2 \cdot 1} \quad \text{ただし、} {}_n C_0 = 1 \text{ とする。}$$

${}_n C_r$  の性質 
$${}_n C_r = {}_n C_{n-r}$$

② 試行と事象

同じ条件のもとで繰り返すことができ、その結果が偶然によって決まる実験や観測を **試行** という。また、試行の結果として起こる事柄を **事象** という。

③ 反復試行の確率

1回の試行で事象  $A$  の起こる確率を  $p$  とする。この試行を  $n$  回繰り返すとき、事象  $A$  がちょうど  $r$  回起こる確率は 
$${}_n C_r p^r (1-p)^{n-r}$$
 ただし、正の数  $a$  に対して、 $a^0 = 1$  と定める。