

# ヒストグラム、散布図

## 1 集団の位置尺度

集団の中心位置を表す尺度

**平均** もっとも多く使われる算術平均であるが、重心を意味する。数学的に取り扱いやすいし、理論的な側面もきわめて重要であるが、万能というわけではなく、裾の影響に大きく左右される欠点をもつ。 $\bar{X} = \frac{1}{n}\{X_1 + X_2 + \dots + X_n\} = \sum_i X_i/n$  データの平均値 (average, mean value) とは、データの値の総和を個数で割ったものをいい、 $\bar{X}$  (エックスバーとよむ) で表す。 $n$  個のデータの総和は  $k$  個に分けた階級値と度数をつかって表すと、 $X_1 + X_2 + \dots + X_n$  の計算の代わりに  $x_1 f_1 + x_2 f_2 + \dots + x_k f_k$  で総和を求める。

**中位数 (中央値)** 小さいものから大きいものへと、大きさの順に並べて、ちょうど真ん中に位置する値で、小さい方も大きい方もそれぞれ 50% ずつとなる値。比較的頑健性を持ち、常識にもかなっている。ただし、データがたくさんであると、大きさの順に並び替えることに労力を要する。

**モード** 度数分布表で求めた度数の最大となっている値をいう。ヒストグラム (度数多角形) がひとつ山の形を単峰型といい、モードが中心傾向を表す。一つ山で対称なときにはこれら 3 つともほぼ同じ値になるが、集団の山形に歪み (L 字型は正の歪み、逆 L 字型は、負の歪みをもつという) があれば、これらは異なった値となる。

## 2 バラツキの尺度

バラツキを表すための尺度

**分散** 各データと平均からのずれ (差) を平方 (square) した値の平均で、非負の値である。ゼロとなるのは、全てのデータが同じ値、すなわちまったく変動がない場合である。

$$s^2 = s_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - (\bar{X})^2$$

**標準偏差** 分散の単位はもとのデータの 2 乗値となるが、この平方根をとったもので、平均と同じ単位となる。Standard Deviation から SD と約されることが多い。

$$s = \sqrt{s^2} = \sqrt{s_X^2}$$

**四分位数** 中位数 (メディアン) がちょうど 1/2 となる点であったと同様に、4 つの 25% ずつに分けた値で、第 1 四分位数 (小さい方から 25% の点)、第 2 (中位数)、第 3 四分位数 (75% の点) といい、第 3 から第 2 を引いた値を 4 分位範囲という。

**パーセンタイル (百分位数)** ; 非常に大きな集団で、これを 100 個のパーセントに分けたものをいう。上位 5% 点、下位 5% 点がとくに集団の特質や特徴を表すために用いられる。

**問 2.1**

5 個の数値データ 3, 5, 4, 1, 7 の平均と分散を計算し、これをもちいて、つぎの場合の数値データの平均と分散を計算しなさい。

(a) 13, 15, 14, 11, 17

(b) 2.3, 2.5, 2.4, 2.1, 2.7

**問 2.2**

4 個の数の平均は 5、分散は 2 で、もうひとつの集団の 6 個の数値の平均は 7、分散は 3 であるという。この 2 つを合わせた 10 個に対する平均と分散を計算しなさい。

**問 2.3**

与えられたデータ  $X_1, \dots, X_n$  について、これを変換して平均は 50、分散は 100、(標準偏差は 10) にするにはどうしたら、よいか。これは「tスコア (偏差値)」とよべます。