

「平均に」強くなろう.....

右を見ても左を見ても、この世は数字だらけです。野球を見れば9対8, オリンピックは10.0, 昼飯食えば400円。こんなに数字だらけの現代に生き、さらに統計を扱うとなると少しは数字の性質や扱い方を知っていた方が便利な時があります。と言ってもむずかしいことは書く方も困難ですから、書く方も読む方も困難をきたさないように進むことにします。統計の数字の基本は、和と平均と比率です。シリーズ「統計を考える」今月は平均を追ってみましょう。

集団を表す代表値には、平均の他に**モード**(mode)と**メジアン**(median)というのが有ります。この3つが代表値の3羽ガラスですが、3羽のうち最も多く飛びまわっているのが御存知「平均ガラス」です。平均にも2種あって、**算術平均ガラス**と**幾何平均ガラス**がいます。この2種を比べてみると前者の方が後者よりも大きく、また人によく好かれます。しかし算術ガラスが人に好かれるのは大きいからではなく扱いやすいからなのです。幾何ガラスは性格がこむずかしくて扱いにくいために嫌われてしまうのです。算術平均の扱いやすい点は何よりもまず第1に計算が簡単なことです。たとえば「2・3・4・5」の算術平均は

$$\frac{2+3+4+5}{4} = 3.5$$

です。これを一般に、

$$\bar{x} = \frac{\sum x_i}{n}$$

で表します。扱いやすい2番目の点は先の式からもわかるとうり、平均値 \bar{x} に項数 n をかければ総和が得られるために利用価値が非常に大きいことです。一方幾何平均というのはグループの n 個の値を全部掛け合わせてそれを n で開いた値です。たとえばここに8人の人がいてその給料が、1万円が2人、2万円が5人、16万円が1人だったとします。この人達の「平均給料」は算術平均では3.5万円ですが、これが幾何平均では

$$\sqrt[8]{1 \times 1 \times 2 \times 2 \times 2 \times 2 \times 2 \times 16} = \sqrt[8]{572}$$

となり平均2.2万円です。実際に皆さん計算してみてください。算術平均を出すように簡単にできましたでしょうか。平均給料は3.5万円より2.2万円の方がおだやかです。1人16万円の高額者のために算術平均がひき上げられてしまったのです。

3羽ガラスのうち1羽は以上のような性質ですが、他の2羽についてもちょっとその姿を観察してみます。

メジアンは和名で中央値と呼ばれるように、ある集団のまん中に位置する値です。ひとつの集団を構成している連中を小さい順から大きい順に並べてその中央にきたのをメ

ジアンというのですから計算も何もありません。「5・5・6・9・9」のメジアンは「6」, 「1・3・6・6・6」のメジアンも中央の「6」です。

モードは和名で最頻値又は並み値と呼ばれます。最も頻ばんに出て来る値だからです。たとえば、

「2・3・3・3・5」のモードは「3」, 「1・3・6・6・6」のモードは「6」です。

さて、これで3羽ガラスの説明は終わりですが、集団を表すこれら3つの代表値にもそれぞれ欠点があります。算術平均は確かにその性質が優れていますが、誰もが気楽に利用し、またどこにでも顔を出して「私がこの集団の唯一の代表者であります。」というデカイ顔をしているので欠点が隠れてしまい、その集団の姿を誤解させたり、また、算術平均値と自分のナマ身を比べて我と我が身を嘆いたり喜んだりする人を絶やしません。それは、その代表値の性質をよく知らないからです。また、そのグループの代表値に何を使うかよく吟味しないで何でもかんでも算術平均で表してしまうからです。ミスもクソも一緒では困るのです。代表者を選ぶときには十分注意しないと後で大変なのはどこの世界でも同じです。オリンピックの体操競技やスキーのジャンプ競技の飛型点の採点方法は、何人かの審判員が出したそれぞれの点数の一番高い点数と一番低い点数を除いた残りの点数の算術平均になっているようです。スポーツの世界でもこのように注意しているのですから、統計の世界ではなおさら注意しなければなりません。

さて、それでは次に3つの集団を見て下さい。

(ア) 7・7・7・7・7

(イ) 5・6・7・8・9

(ウ) 1・4・7・10・13

上の3つのグループの特徴を見るために先の3羽ガラスに登してもらおうと、平均もメジアンも(ア)~(ウ)それぞれ同じく7です。どの集団も7で代表されるとはいえ、明らかにこの集団の間には差異が有るのですが、その差異が3羽ガラスでは明らかにならないのです。この場合代表値としての3羽は失格しました。どこが違うのでしょうか?違うのはバラツキなのです。ではそのバラツキ方を調べて表現するにはどんな方法が有るのでしょうか。それにはまず**レンジ**があります。レンジはRangeで範囲のことです。

(ア)の場合は7-7=0

(イ)の場合は9-5=4

(ウ)の場合は13-1=12

富永重己

です。計算といえはひき算だけのすこぶる簡単な方法です。しかし計算が簡単な分、いかにも荒っぽいやり方だという感じがします。たとえば、

- (イ) $1 \cdot 1 \cdot 1 \cdot 2 \cdot 9$
- (オ) $1 \cdot 3 \cdot 5 \cdot 7 \cdot 9$

の2つのグループはバラツキがともに「8」のレンジで表わされますが、この2つのグループのバラツキ方には明らかに差異が有ると見るのが普通です。そのバラツキ方の差異がレンジで表しきれないといえます。レンジは簡単にらせてある程度バラツキを表現できて便利ですが、その計算はほんのわずかな値だけ取り上げるだけで、その他大勢のデータを無視している分、荒っぽいのです。これは前の3羽ガラスにもいえることです。算術平均が代表値としてモードやメジアンよりも優れている点は、算術平均が集団の中のデータすべてを計算に入れていることなのです。それと同じように、バラツキを調べる上でもデータをすべて考慮して表現する良い方法があればこれは便利です。そしてそんな方法が実際有るのです。それが「標準偏差」なのです。どうして標準偏差と呼ばれるのかわかりませんがそう呼ばれています。これは重宝です。

標準偏差はシグマで表されますが、このシグマは大文字の Σ ではなく（御存知のとおり大文字は総和を表します）小文字の「 σ 」です。このやり方は、ある集団の算術平均値（ \bar{x} ）から個々の単位（変量 x ）がどのくらい離れた距離に有るか調べる方法です。つまり「 $x - \bar{x}$ 」です。個々の単位についてですから、ひとつひとつ「 $x - \bar{x}$ 」を出してそれらすべてを合計します。つまり

$$\sum(x - \bar{x})$$

で表されます。統計学ではこの \bar{x} からの x の距離を偏差と呼んでいるのです。しかしちょっと待って下さい。実は、算術平均の性質を考えて頂ければおわかりのとおり、 $\sum(x - \bar{x})$ はいつの場合でも「0」になってしまうのです。例えば先の(オ)の例でやってみますと、

$1 \cdot 3 \cdot 5 \cdot 7 \cdot 9$ の平均 \bar{x} は5ですから $1 - 5 = -4$ 、 $3 - 5 = -2$ 、 $5 - 5 = 0$ 、 $7 - 5 = 2$ 、 $9 - 5 = 4$ これを合計すると、 $(-4) + (-2) + (0) + (2) + (4) = 0$ です。そこで考えられる方法としては、 $+$ ・ $-$ の符号を取りはずして絶対値で計算するやり方で、上の例では12になります。これを x の数 n で割って \bar{x} からの平均的な距離を出すと2.4でこれならまあ良さそうですが、これは実は平均偏差と呼ばれているものなのです。標準偏差に行き着くた

めにはもうひと工夫必要です。 $+$ ・ $-$ をはずす方法は、絶対値を使うよりも後々の計算のことを考えてもっと優れた方法を使った方が良いのです。それは2乗して、一の符号を+に変えてしまうやり方です。(イ)の例を使って早速やってみましょう。式は

$$\frac{\sum(x - \bar{x})^2}{n}$$

となるわけです。

$(-2)^2$ は4、 $(-4)^2$ は16になりますから $\sum(x - \bar{x})^2$ は40です。これを n で割って平均的な距離を出すと「8」と出ますが、実はこれは2乗して計算したため σ^2 の値になっています。 σ を出すには $\sqrt{\quad}$ で開いて元に戻してやらなければなりません。

$$\sqrt{8} = 2.83, \text{これが}\sigma\text{です。}$$

こうして、標準偏差を出す式は

$$\sigma = \sqrt{\frac{\sum(x - \bar{x})^2}{n}} \quad \text{..... ①}$$

と決めます。この式は覚えておいた方が良いでしょう。実際の例で標準偏差を出す場合には平均値 \bar{x} が小数点以下の細かい数字になることが多く、いちいち $(x - \bar{x})^2$ を出さなければならない①式のやり方では手間がかかります。そのため、もっと簡便なやり方として①式を変形した公式、

$$\sigma = \sqrt{\frac{\sum x^2}{n} - \bar{x}^2} \quad \text{..... ②}$$

を使います。御用とお急ぎでない方でもこの方が便利です。

σ の値は0に近いほど、分布の広がり狭いことを示します。先の(イ)の例を使って σ の値を出してみて下さい。

x (イ)	x^2	$x - \bar{x}$	$(x - \bar{x})^2$
1	1	-1.8	3.24
1	1	-1.8	3.24
1	1	-1.8	3.24
2	4	-0.8	0.64
9	81	6.2	38.44

$$n = 5, \quad \sum x^2 = 88, \quad \sum(x - \bar{x})^2 = 48.80$$

$$\bar{x} = 2.8 \quad \bar{x}^2 = 7.84$$

①のやり方では、

$$\sigma = \sqrt{\frac{\sum(x - \bar{x})^2}{n}} = \sqrt{\frac{48.80}{5}} = \sqrt{9.76} \approx 3.12$$

②のやり方では

$$\sigma = \sqrt{\frac{\sum x^2}{n} - \bar{x}^2} = \sqrt{\frac{88}{5} - 7.84} = \sqrt{9.76} \approx 3.12$$

(イ)の標準偏差は2.83でしたので(オ)の方が(イ)の集団より分布の広がり狭いことがわかりました。（県消費統計係）