

# 百均の電卓で解ける 統計入門

数学的な厳密さよりも直感的に理解することに  
重点を置いています・・・

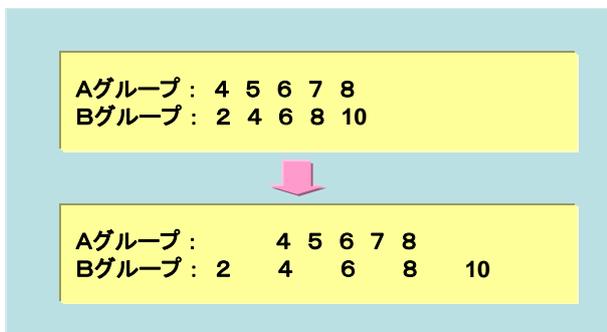
## 散らばり? - I -

前回までに、データが数値的に広がって存在していて、その広がりの中から「全データを代表する数値」として**平均値**のお話をしてきました。

データたちは、平均値の周辺にどのように分布しているかを見るのが**ヒストグラム**でしたね。

しかし、平均値やヒストグラムでは、データがどのくらい広がっているのか、あるいは散らばっているのかはわからないのです。

で、そのデータたちが、どのくらい広がっているのか、ばらついているのかを考えていきましょう。



### (1) 平均値からの距離?

皆さん、第 1 回で「図の 2 つのグループを比べると、どちらがどのくらいバラツキが大きいと思いますか? 並べ変えると・・・多くの方が『BはAより 2 倍ぐらいばらついているかなあー』と感じるのではないのでしょうか。」というお話をしたことを覚えていますか?

これを例題に考えていきましょう。

前号までに平均値の求め方はやりましたので、2 つのグループはともに平均値が「6」であることはすぐに理解できると思います。

では、各データが平均値からどのくらい大きいのか、あるいは小さいかを見てみましょう。

A グループ	-2	-1	0	1	2
B グループ	-4	-2	0	2	4

$$A \text{ グループ} : \{(-2)+(-1)+(0)+(1)+(2)\} \div 5 = 0 \div 5 = 0$$

$$B \text{ グループ} : \{(-4)+(-2)+(0)+(2)+(4)\} \div 5 = 0 \div 5 = 0$$

このように、平均値より大きいときはプラス、小さいときはマイナスで表した数値を統計学では『**偏差(deviation)**』と呼びます。私たちが知りたいのは、この 5 つの偏差を 1 つの数字で代表させたものですね。しかし、単純に算術平均(足して個数で割ったもの)を求めると・・・

実は、どんなデータに関しても、単純に偏差の算術平均をとると「0」になってしまうのです。ここでは、その証明は割愛しますが、算術平均を使う方法がうまくいかないことは、直感的にわかってもらえると思います。

では、どうしたらよいのでしょうか?

### (2) 符号をなくす?

こういった、不具合が起こる原因が符号であることに気付きましたか? 平均値より大きいときはプラス、小さいときはマイナスで表した数値をそのまま使ったので打ち消し合いが起きてしまったのです。

そこで、マイナス側の偏差とプラス側の偏差を同等に扱うためには、絶対値にするかわりに二乗するのです。つまり、偏差の二乗和の平均をとるわけです。

$$A \text{ グループ二乗和の平均} = \frac{(-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2}{5} = 2$$

$$B \text{ グループ二乗和の平均} = \frac{(-4)^2 + (-2)^2 + (0)^2 + (2)^2 + (4)^2}{5} = \frac{40}{5}$$

ここで出てきた数値量を『**分散(variance)**』と呼びます。

ただ・・・これでは、ここの偏差より大きな数字となってしまう、「全データを代表する数値」とは言いにくいですね。

また、この例では単位がありませんが、m とか、分であったら、単位も二乗になってしまいます。cm であったら、面積になってしまいますね。

そこで、第 2 回で述べた「**二乗平均**」を思い出してください。それぞれの差を二乗して足し、個数で割り、さらにルートする方法です。



$$A \text{ グループ二乗平均} = \sqrt{\frac{(-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2}{5}} = \sqrt{2}$$

$$B \text{ グループ二乗平均} = \sqrt{\frac{(-4)^2 + (-2)^2 + (0)^2 + (2)^2 + (4)^2}{5}} = \sqrt{\frac{40}{5}}$$

この統計量を「**標準偏差(Standard Deviation)**」と呼びます。この英語の頭文字を取って、「**S.D.**」と略されるのです。

この項 ⇒