

雑感 ヒストグラムと GeoGebra

■ 「データの分析」で箱ひげ図や散布図を描くとき、GeoGebraの「表計算・統計」が便利なのでよく利用する。

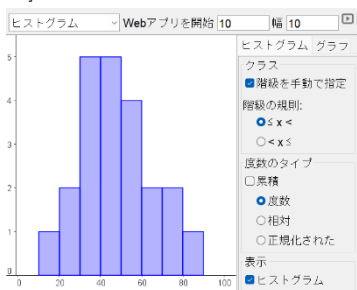
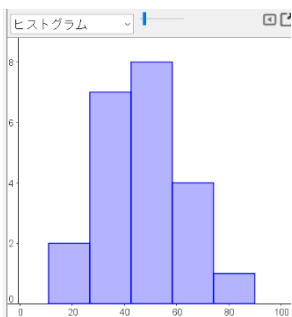
「1変数解析」の中にヒストグラムもあるのだが、使い勝手がイマイチのように感じて使うことが少なかった。

例えば、{11, 25, 28, 31, 33, 37, 36, 38, 41, 44, 47, 45, 49, 50, 53, 54, 55, 65, 66, 70, 71, 90} という22個のデータをヒストグラムにすると、最初、右のような階級数が5

のものが現れる。[この最初の階級数5は、スタージェスの公式※からの $1 + \log_2 22 \approx 5.45$ によるものかも知れない]

この階級数は、上の青い | を動かすことによって変えることができるが、階級の区切りがどうなっていくのかよく分からないので、実用的ではないように思われる。

しかし、☑ をクリックして階級条件を指定すれば、右のように希望の設定が可能になっている(恥ずかしいことに、先日までこれを知らなかった)。



■ ところが、表示されたヒストグラムを子細に見ると、データの内、90が正しく描かれていない。設定の階級の規則から、本来90は階級 $90 \leq x < 100$ に長方形が描かれるべきデータであるが、階級 $80 \leq x < 90$ に描かれてしまっている。

[この90を91などに変更してみると、なぜか正しく描いてくれる] このエラーと思われる現象は、Yahoo知恵袋で質問の形で指摘されたものである。

回答で参照サイト https://wiki.geogebra.org/en/Histogram_Command を紹介したところ、質問者からは「一番右側の階級は「 $a \leq x \leq b$ 」となるよう」だとの返答があった。実際確認してみると

By convention this uses the $a \leq x < b$ rule for each class except for the last class which is $a \leq x \leq b$

とある。これは、100点満点の試験のとき、満点の100点を90点台の階級に含めてしまうような流儀なのであろう。90でなく91だと正しいヒストグラムになるのは、上の注記では不十分な理由が何かあるのだろう。

■ そういった意味では、こういった微妙なルールを知っていないと間違えてしまうことになり、要注意である。

では、GeoGebraではこの「正しい」ヒストグラムは描けないのかと言うと、面倒を覚悟すればコマンド

```
Histogram( {10, 20, 30, 40, 50, 60, 70, 80, 90}, {11, 25, 28, 31, 33, 37, 36, 38, 41, 44, 47, 45, 49, 50, 53, 54, 55, 65, 66, 70, 71, 90}, false)
```

で可能であり、グラフィック画面に右のように正しく描かれる。

2つ目の正しくないヒストグラムは、どうやら

```
Histogram( {10, 20, 30, 40, 50, 60, 70, 80, 90}, {11, 25, 28, 31, 33, 37, 36, 38, 41, 44, 47, 45, 49, 50, 53, 54, 55, 65, 66, 70, 71, 90}, true, 10)
```

といったコマンドで描かれるように仕組みられているようである(右図)。

■ 奥が深く、知らないことが多い。

※ 最近、「スタージェスの経験則」といった記述を見かけるが、彼の名誉のために言うておけば、経験則ではないはず ([hist.pdf\(plala.or.jp\)](http://hist.pdf(plala.or.jp)))。