

# 数学 I 『データの分析』について

～四分位数、箱ひげ図など～

竹田 裕一<sup>1</sup>

新聞・雑誌やインターネット上に様々な指標や調査結果（株価指数や世論調査など）があるが、その結果を正しく理解するためには統計の知識が必要となる。このため新学習指導要領では、小学校・中学校・高等学校に資料の活用（統計）が必修として入ることになった。

今回の講演会では数学 I の「資料の活用：データの分析」の内容のうち、従来の数学 B・C に含まれない「四分位数」、「四分位偏差」、「箱ひげ図」を中心に解説を行う。

## 1 中学で新しく教える範囲

現在の学習指導要領において数学 B および数学 C で学習する内容の一部が中学の数学で学習することになった。数学 I で新しく教える「四分位数」などは、この内容を踏まえているので復習を兼ねて定義などを簡単に解説しておく。

### 1.1 分布を代表する統計量と散らばり具合を表す統計量（数学 B ⇒ 中学 1 年）

分布を代表する統計量として 3 つの代表値「平均 (mean)」、「中央値 (median)」、「最頻値 (mode)」がある。これらの違いを簡単に解説すると

#### ・平均値 (mean)

数値的な中心を表す代表値で、すべてのデータの値  $(x_1, x_2, \dots, x_n)$  を足しその個数  $(n)$  で割った値である。

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

#### ・中央値 (median)

確率的な中心を表す代表値で、大きさの順に並べたとき中央にくる値である。偶数個の場合は、中央に並ぶ値の平均値を中央値とする。

#### ・最頻値 (mode)

頻度に関する代表値で、頻度が最大（最もよく現れる）となる値である。通常は度数分布表に整理したときに、最大度数をもつ階級値を使う。

また平均値からの散らばり具合を表す統計量として「分散 (variance)」がある。一般に分散の値が小さいほど平均値の近くにデータがあることを示している。各データの値  $x_i$  の値と平均値  $\bar{x}$  の差を  $x_i$  の偏差とよび、分散 ( $s^2$ ) と標準偏差 ( $s$ ) の定義は次のとおり。

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$
$$s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

<sup>1</sup>神奈川工科大学 基礎・教養教育センター  
E-mail : y-takeda@ctr.kanagawa-it.ac.jp

## 1.2 母集団と標本（数学C ⇒ 中学3年）

調査対象の性質・特徴などを調査するときに、調査対象すべて（**母集団**）を調べることが可能であればすべてのデータを調査する（**全数調査**）ことによって正確な結果を得ることができる。しかしながら、全数調査が物理的な問題（無限個ある場合や破壊検査など）や時間・費用の問題で実施できないことが多い。その場合、母集団の一部分（**標本**）を調査する（**標本調査**）ことによって、性質・特徴を推測することになる。正しい推測をするためには、母集団を正しく定義することと、母集団から偏りなく標本を取り出すことが重要となる。

例. 選挙

選挙は1票でも多い人が当選するので必ず全数調査が行われている。しかしながら、選挙速報の「当選確実」は出口調査などの標本調査から推測された結果である。推測であるために、当選確実を間違えることや、票差が僅差であるとき、当選確実が出ないことがある。

調査事項：誰が当選するか

母集団：投票した人すべて

標本：ある投票所で投票した人の一部

（いわゆる出口調査）

選挙の場合、母集団を有権者すべてと考える場合もあるが、当選者を決めるのは投票数であるから投票した人のみを母集団と考えればよい。

前節で解説した代表値は、母集団の真の値（**母平均**や**母分散**）ではなく標本を使って推定した値（**標本平均**や**標本分散**）である。推定値の計算方法はいろいろな種類があり、それぞれ異なる基準に従って作られたものである。母平均  $\mu$  の推定量である標本平均  $\bar{x}$  は前節で定義した式が使われることが多いが、母分散  $\sigma^2$  の推定量である標本分散  $s^2$  は多くの本で2つの推定量が示されている。

不偏推定量（数値的に近い）

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

最尤推定量（確率的にあてはまりが良い）

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

一般に母分散  $\sigma^2$  の推定量は数値的に近い不偏推定量が使われることが多い。Excel で分散の値を求める場合、不偏推定量を VAR、最尤推定量を VARP で求めることができる（正確には、関数 VARP は与えられたデータが母集団すべてである場合の母分散の値を求めている）。

## 2 高等学校で新しく教える範囲

新学習指導要領において、数学 I において学習することになった「データの分析」の内容は、現在の数学 B の「統計とコンピュータ」をほぼそのまま移行した形になっている。現在の内容からの変更点は、いくつかの用語の変更（資料 ⇒ データ、相関図 ⇒ 散布図）と、代表値の内容が中学に移ったことから、新しい内容として「**四分位数**」, 「**四分位範囲**」, 「**箱ひげ図**」が加えられた点にある。「四分位数」や「四分位範囲」は分散や標準偏差と同様に分布の散らばり具合を表す統計量であり、「箱ひげ図」は「四分位数」を使って、分布の散らばり具合をグラフ化し直感的にわかりやすくしたものである。次節以降で「四分位数」, 「四分位範囲」, 「箱ひげ図」の解説と、平均や分散との違いについて述べる。

## 2.1 四分位数と四分位範囲

**四分位数**はデータを大きさの順に並べたとき、4つに分割した所にある3つの数であり、値の小さなほうから**第1四分位数**、**第2四分位数**、**第3四分位数**とよぶ。第2四分位数についてはデータを大きさの順に並べて中央に来る値であるから中央値と一致している。また、データの散らばり具合を第1四分位数、第2四分位数、第3四分位数と最大値、最小値の5点で表す方法があり、**5数要約**とよぶ。ちなみに、四分位数と同じ考え方で**三分位数**や**五分位数**なども定義可能である。

### 2.1.1 四分位数の求め方

四分位数は次のように求める。

- 1) データを大きさの順に並べ、中央値を求め、第2四分位数とする。
- 2) 中央値のデータを含まない中央値以下のデータのみで中央値を求め、第1四分位数とする。
- 3) 中央値のデータを含まない中央値以上のデータのみで中央値を求め、第3四分位数とする。

データの個数が偶数の場合は、中央値のときと同様に中央に並ぶ値の平均値を中央値とする。

例1. 次のような15個のデータの場合

順位	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
データ	11	11	12	16	16	16	18	18	19	20	20	20	22	23	24

- 1) 第2四分位数(中央値)を求める

8番目が中央の値なので第2四分位数は18

- 2) 第1四分位数を求める

1～7番目のデータの中央値を求めればよい。

4番目が中央の値なので第1四分位数は16

- 3) 第3四分位数を求める

9～15番目のデータの中央値を求めればよい。

12番目が中央の値なので第3四分位数は20

最小値	第1四分位数	第2四分位数	第3四分位数	最大値
11	16	18	20	24

例2. 次のような8個のデータの場合

順位	1	2	3	4	5	6	7	8
データ	10	11	11	13	14	18	20	21

- 1) 第2四分位数(中央値)を求める

4番目と5番目が中央の値なので第2四分位数は  $\frac{13+14}{2} = 13.5$

- 2) 第1四分位数を求める

1～4番目のデータの中央値を求めればよい。

2番目と3番目が中央の値なので第1四分位数は  $\frac{11+11}{2} = 11$

- 3) 第3四分位数を求める

5～8番目のデータの中央値を求めればよい。

6番目と7番目が中央の値なので第3四分位数は  $\frac{18+20}{2} = 19$

最小値	第1四分位数	第2四分位数	第3四分位数	最大値
10	11	13.5	19	21

### 2.1.2 四分位範囲の求め方

四分位範囲はデータの中心付近の50%が含まれる区間の大きさを示す。

$$\text{四分位範囲} = \text{第3四分位数} - \text{第1四分位数}$$

また、四分位範囲の半分の大きさを四分位偏差とよぶ。

$$\text{四分位偏差} = \text{四分位範囲} \div 2$$

この他、データの100%が含まれる区間の大きさを示す範囲がある。

$$\text{範囲} = \text{最大値} - \text{最小値}$$

実際に前節の例の場合のそれぞれの値を求めると

例1. の場合

範囲	四分位範囲	四分位偏差
$24 - 11 = 13$	$20 - 16 = 4$	$4 \div 2 = 2$

例2. の場合

範囲	四分位範囲	四分位偏差
$21 - 10 = 11$	$19 - 11 = 8$	$8 \div 2 = 4$

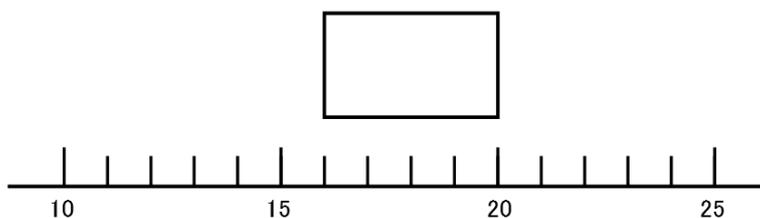
ここで定義した四分位数や四分位範囲はあくまでも真の値を推定した値であり、1.2節でも述べたように推定値の計算方法はいろいろな種類があり、それぞれ値が異なる。そのため、ソフトによってはここで定義した値と異なる値が出てくることがある。実際 Excel で四分位数を求める関数 QUARTILE を使って四分位数を求めると、ほとんど同じ値が得られるが、例2の第3四分位数のみ19ではなく18.5が得られる。

## 2.2 箱ひげ図のかき方

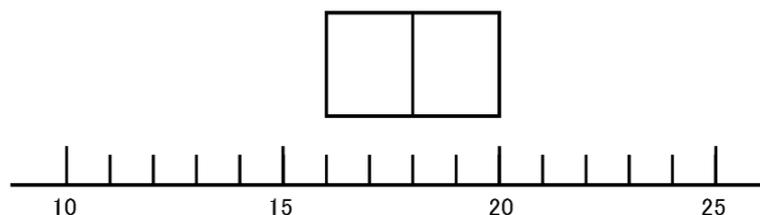
箱ひげ図は5数要約を元に次のようにしてかく。

例1の場合

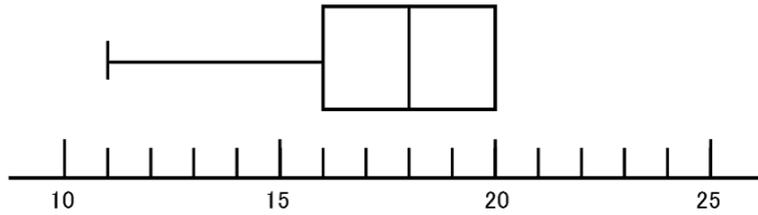
- 1) 第1四分位数と第3四分位数を両端とした箱（長方形）をかく。  
(長方形の幅が四分位範囲に対応する)



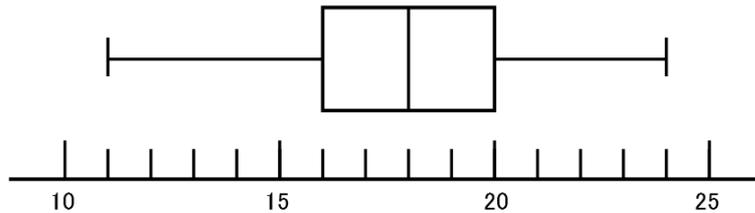
- 2) 第2四分位数（中央値）に線を引く。



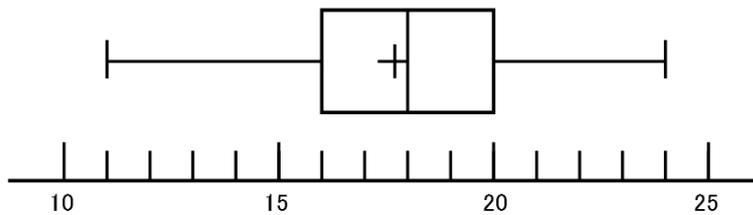
3) 最小値に線を引き、箱の左側までひげ(線)をかく。



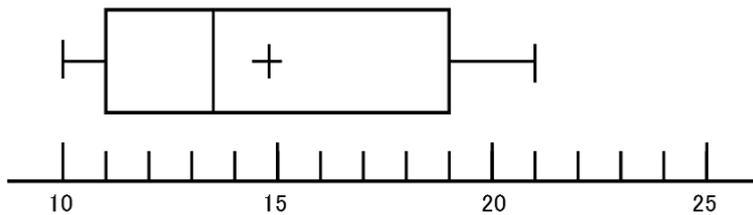
4) 最大値に線を引き、箱の右側までひげ(線)をかく。



一般に箱ひげ図は4)のような形になるが、場合によっては平均値(17.73)に印をつけることがある。



例2の場合の箱ひげ図をかくと次のようになる。



上記の箱ひげ図は横方向であるが、株価のローソク足のような縦方向の箱ひげ図をかくこともある。

箱ひげ図を見ると、「データが多く集まっている所」や「分布が左右対称」などがわかる。実際に例1はほぼ左右対称になっていて、中央値を含む50%のデータは中央値の近くにまとまっていることがわかる。例2は前半50%のデータは10から13.5の狭い範囲にあるが、後半50%のデータが13.5から21と範囲の幅が倍以上になっていることから、明らかに左右対称ではない(右側のすそが長い分布)ことがわかる。このように箱ひげ図を見るとある程度ヒストグラムの形も予測できる。

**演習.** 次のデータの5数要約を求め、箱ひげ図を完成させよ。

番号	1	2	3	4	5	6	7	8	9	10	11	12	13
データ	21	20	22	10	32	18	11	13	15	20	24	16	25

(平均値  $\bar{x} = 19$ )

### 3 まとめ

新聞やインターネットに載っている情報の多くは、分布の真ん中を示す値として平均値のみが書かれていることが多く、中央値が書かれていることが少ない。これは母集団分布が左右対称であれば「平均値＝中央値」が成立するために平均値を載せておけば十分だからである。特に母集団分布が正規分布の場合、母平均  $\mu$  と母分散  $\sigma^2$  がわかれば分布の形が決まってしまうため、平均と分散さえ推定すれば分布の形も推定できる ( $[\mu - \sigma, \mu + \sigma]$  の間に 68%,  $[\mu - 2\sigma, \mu + 2\sigma]$  の間に 95% のデータが入る)。

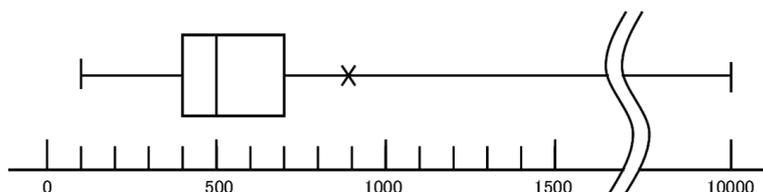
母集団分布が左右対称で無い場合、平均値は直感的な真ん中の値を示していないことがある。例えば所得の例のように一部の人間が大きな収入がある場合、平均値は感覚よりも大きな値となってしまう。この場合平均値よりも中央値の方が感覚と一致している。

例. 20人の年収が次のようなとき、平均値・中央値はいくらになるか。

100万:	1名	200万:	1名	300万:	2名	400万:	3名
500万:	4名	600万:	3名	700万:	3名	800万:	1名
1000万:	1名	1億:	1名				

答) 平均値: 990万円      中央値: 500万円

上記のデータでは、9割の人が平均値以下の年収となってしまう真ん中の値とはいえない。これは明らかに年収1億円の人がいるために、平均値が上がってしまったためである。このデータに対して箱ひげ図をかくと次のようになる (実際に箱ひげ図をかくときは1億円を外れ値として扱うこともある)。



(第1四分位数: 400万円, 第3四分位数: 700万円)

箱ひげ図を見ると平均値が箱の中に入っておらず、明らかに右すそが長い分布であることがわかる。実際の年収の分布は厚生労働省が毎年公表しており、「平成20年 国民生活基礎調査の概況」

(<http://www.mhlw.go.jp/toukei/saikin/hw/k-tyosa/k-tyosa08/index.html>) をみると平均値556万円に対して中央値448万円と100万円以上の差がある。この他に五分位数が公表されており、

第1: 210万円      第2: 361万円      第3: 549万円      第4: 814万円

である。詳細なデータが無いと最大値が大きすぎると考えられるので箱ひげ図をかくことができないが、実際のデータも右すそが長い分布であることがわかる (実際にヒストグラムが載っており、右すそが長い)。そのため平均値のみが一人歩きすると、自分の年収が平均収入よりも低いと嘆く人が多いことになる (平均年収より100万円以上少ない人が50%以上いる)。

この例のように母集団分布の左右対称が崩れた場合、平均・分散を使うより中央値や四分位数を使ったほうが適当である。さらに左右対称であるかどうかを判断するためにはヒストグラムや箱ひげ図などのグラフを見ると傾向がわかる。ヒストグラムは階級の幅を変えると、分布の形も変わってしまうことがあるが、箱ひげ図は一意に決まる5数要約を使ってかくので、必ず同じ形となる利点がある。

現在のデータ分析はコンピュータにデータを入力すればいろいろな指標 (統計量) やグラフを簡単に得ることができるが、それ故にあきらかに間違った結果や、ミスリードを誘うような結果も散見される。小学校から大学までうまく連携して正しい統計の知識を教えることによって、目先の数字にだまされずに正しい判断ができるような人を育てて行きたいと思っています。