

「統計じかけのオレンジ」

— A Statistics-work Orange —

第1回 平均、分散、標準偏差

一般社団法人 日本下水道施設業協会

技術部長 堅田 智 洋



1. はじめに

統計学は記述統計学と推測統計学の2つに大きく分類することができる。

記述統計学は、調査や実験等で集めたデータを、整理、表やグラフ化、数値化することでイメージ的、直感的に理解しようとするものである。一方、推測統計学は、統計学の手法と確率理論を使って、「全体の把握が難しいくらい大きな対象」や「まだ起きていない、未来に起きること」といった「全体像」を、その時点で得られている「部分的な結果」から推測しようとするものである。

今回、筆者は、推測統計学のほんのさわりの部分を理解できた範囲でまとめてみた。ただ、不正確な記述も散見されると思われるため、実務への活用の際は、他の専門書で正確な情報を確認してください。

2. さまざまな統計量

2.1 統計量とは

統計量とは、記述統計学において、あるデータの特徴を1つの数値に要約して表現したものであり、よく知られているものに、平均値、分散、標準偏差、中央値、最頻値、最大値、最小値がある。ここでは、その中でも推測統計学においても重要となる平均値、分散、標準偏差を概説する。

2.2 平均値

平均値は「データを合計したものをデータ数で割った値」である。これは「算術平均」と呼ばれるもので、例えば、2つの数 x と y の算術平均は $\frac{x+y}{2}$ で求まる。

平均値の意味合い、性質として重要なポイントを以下に挙げる。

- ・データをヒストグラムで表すと、平均値はヒストグラムをやじろべえとみなした時のつり合いの支点となる。
- ・データは平均値の周辺に分布している。
- ・多く現れるデータの平均値への影響力が大きい。

※ヒストグラム：データの分布状況を視覚的に表現するため、測定値の存在する範囲をいくつかの区分に分け、各区分を底辺としてその区間に属する測定値の出現（相対）度数に比例する面積をもつ長方形を並べた図。

実は、平均の求め方はこの「算術平均」のほか「相乗平均（あるいは幾何平均）」、「二乗平均」、「調和平均」等がある。「二乗平均」は各データを2乗して合計し個数で割ってその後ルートにしたもので、2つの数 x と y の二乗平均は $\sqrt{\frac{x^2+y^2}{2}}$ である。この「二乗平均」の手順は後で標準偏差のところでも出てくるので、算術平均と対比して示しておく（図2-1）。

二乗平均は、各データを最初に二乗してから算術平均を行い、最後に、最初の二乗操作を戻すためにルートを行うのだ、と解釈するとわかりやすいかもしれない。

他の平均の説明は割愛するが、いずれも「 x と y の間にある1つの数」を代表値として選び出していることは共通している。「 x と y を1つの数で代表するのにどの平均がふさわしいのか」は「そのデータ全体に関して何を知りたいのか」に依存して決まるため、用途に従って使い分けることになる。

本稿では、以降、特に断りがない限り「平均」は算術平均を指すこととする。

STEP	算術平均	二乗平均
1	各データ x, y	各データ x, y
2	↓	各データを二乗する。 x^2, y^2
3	各データを合計する。 $x+y$	「各データの二乗」を合計する。 x^2+y^2
4	各データの合計を データ数で割る。 $\frac{x+y}{2}$	「各データの二乗」の合計を データ数で割る。 $\frac{x^2+y^2}{2}$
5	—	「『各データの二乗』の合計を データ数で割ったもの」を ルートする。 $\sqrt{\frac{x^2+y^2}{2}}$

各データの二乗操作

「各データを二乗したもの」の算術平均

全体をルートし、最初の二乗操作をリセット

図2-1 算術平均と二乗平均の手順

2.3 分散、標準偏差

2.3.1 データのバラツキの重要性

私たちは、日常生活においては何かと平均値だけで物事を判断することも少なくないが、実際には平均だけで事象の様子がわかったというわけにはいかないことも多い。例えば、「ある歌手のファン層の平均年齢は25歳である。」という話を聞いたとする。それだけで、この歌手のファンには20歳代の若者が多いのだろうと早とちりしてはいけない。子供から高齢者まで幅広い年齢層から支持されていることが、平均年齢という統計量においては25歳という数値で表されたに過ぎない可能性もある。ファンの年齢のばらつきの情報は与えられていないのだから、年齢構成を推し測ることはできないのである。

上記は極端な例だとしても、データの特徴を把握するためには平均値だけでなく散らばりやバラツキを知ることが非常に重要である。小島は自身の著書¹⁾で、そのことを「バスの運行状態」を例にとってわかりやすく説明しているので、以下に要約して紹介する。

ある目的地へ行くために、バスAとバスBのどちらかを利用しようとしている。バスAは、これ

から乗車しようとするバス停に、予定到着時刻に対して等確率で2分遅れたり2分早く来たりする。同じく、バスBは、等確率で10分遅れたり10分早く来たりする。この場合、バス停への到着時刻を平均値だけで見る分には、どちらのバスも時刻表通りに運行しているバスと見なすことができ、甲乙つけられない。しかし、実際の到着時刻が予定時刻からどの程度前後するか（ばらつくか）を考慮すると、実際には、バスAを選択するのではないだろうか。上記のバスAにおける「2分」とバスBにおける「10分」がダイヤの乱れ、つまりバス停到着時刻のバラツキ、散らばり具合を表している統計量だと考えることができる。そして、利用するバスの選択には、平均値よりこの散らばり具合を知ることの方が重要だということになる。

2.3.2 分散、標準偏差の算出方法

前述したデータの散らばり具合を示す統計量が、分散と標準偏差である。その意味と算出方法を、A高校1年の5人の女子の体重(kg)を例²⁾に以下に示す。

その前にまず、「偏差」という用語がある。偏差とは各データから平均値を引いて得られる値で、

番号	体重 (kg)
1	51
2	49
3	50
4	57
5	43

平均値の算出

$$\frac{51 + 49 + 50 + 57 + 43}{5} = 50 \text{ (kg)}$$

図2-2 平均値の算出

番号	体重 (kg)	偏差 (kg)
1	51	51 - 50 = 1
2	49	49 - 50 = -1
3	50	50 - 50 = 0
4	57	57 - 50 = 7
5	43	43 - 50 = -7

図2-3 各データの偏差の算出

各データが平均値からどれだけ離れているかを表す (図2-2、図2-3)。

偏差 = データの数値 - 平均値 … (式2-1)

私たちは、調査から得た複数のデータを分析する場合、個々のデータの偏差がどのようなものかということよりも、データ全体の傾向としての散らばり具合を知りたいことの方が多いだろう。そこで、図2-3で求めた全データの偏差を平均すればいいと考えるのだが、図2-4で示すように各データの偏差を単純に算術平均するとゼロになってしまう。そもそも、平均値は、平均値より大

きかったり小さかったりする各データを均 (なら) したもので、プラス、マイナスどちらの値も存在する偏差を算術平均すると打ち消しあってゼロになるのは当然である。

そこで、「『偏差を二乗したもの』の平均値」をとってみる。偏差を二乗してから平均を取る (二乗したものを合計し、データ数で割る) ことでプラス・マイナスが打ち消しあわないようにできるので、資料のバラツキの指標となりうる。この値が「分散」である (図2-4)。

番号	体重 (kg)	偏差 (kg)
1	51	51 - 50 = 1
2	49	49 - 50 = -1
3	50	50 - 50 = 0
4	57	57 - 50 = 7
5	43	43 - 50 = -7

偏差の単純な算術平均値は0になってしまう。

$$\frac{1 + (-1) + 0 + 7 + (-7)}{5} = 0 \text{ (kg)}$$



そこで
偏差を二乗してから算術平均を行う。

$$\frac{1^2 + (-1)^2 + 0^2 + 7^2 + (-7)^2}{5} = 20 \text{ (kg}^2\text{)}$$

これが「分散」である。

図2-4 分散の算出

$$\text{分散} = \frac{(\text{偏差の2乗}) \text{の合計}}{\text{データ数}} \quad \dots (\text{式} 2 - 2)$$

しかし、この「分散」には2つの問題がある。ひとつはバラツキを表す数値として大きすぎること、もうひとつは単位が変わってしまっている(kg→kg²) ことである。この2つの問題はいずれも、偏差を二乗してから平均を取るために生ずる。そこで、それを解消するために、最後にこの分散のルートをとる。これが「標準偏差」である。

先述した「二乗平均」の定義から、「偏差の二乗」の平均値のルートをとると、「偏差の二乗平均」になるから、次式が得られる。

$$\text{標準偏差} = \sqrt{\text{分散}} = \text{偏差の二乗平均} \dots (\text{式} 2 - 3)$$

標準偏差は、「各データの偏差の平均値」であり、すなわち、データ全体の「散らばり具合」を示すものである。「各データの偏差（平均値からの離れ方）の二乗平均（二乗し、合計し、データ数で割り、ルートにしたもの）」であるから、「各データの偏差（平均値からの離れ方）の平均」を表しているのである。

例題では、標準偏差 = $\sqrt{20} = 4.472 \approx 4.5$ (kg) となり、各人の体重の平均体重 (50kg) からのバラツキの標準的な幅が4.5kgであるとわかる。

<次号に続く>

【本稿の全体構成】

1. はじめに
2. さまざまな統計量
 - 2.1 統計量とは
 - 2.2 平均値
 - 2.3 分散、標準偏差
 - 2.3.1 データのバラツキの重要性
 - 2.3.2 分散、標準偏差の算出方法
 - 2.3.3 標準偏差の意味
3. 正規分布
 - 3.1 正規分布の特徴
 - 3.2 標準正規分布
4. 推定
 - 4.1 統計的推定とは
 - 4.2 統計的推定のパターン別アプローチ

4.2.1 統計的推定のファーストアプローチ

4.2.2 正規母集団における母分散 σ^2 がわかっているときの母平均 μ の推定

4.2.3 正規母集団における母平均 μ がわかっているときの母分散 σ^2 の推定

4.2.4 正規母集団における母平均 μ がわからないときの母分散 σ^2 の推定

4.2.5 正規母集団における母分散 σ^2 がわからないときの母平均 μ の推定

5. 検定

5.1 検定とは

5.2 検定のパターン別アプローチ

5.2.1 母平均の検定

5.2.2 t検定

5.2.3 カイ二乗検定

【参考文献】

- 1) 完全独習 統計学入門 小島 寛之 2006年9月
ダイヤモンド社
- 2) まずはこの一冊から 意味がわかる統計解析 涌井
貞美 2013年2月 ベレ出版

シリーズ・技術調査報告

「統計じかけのオレンジ」

—A Statistics-work Orange—

第2回 標準偏差、正規分布

一般社団法人 日本下水道施設業協会

技術部長 堅田 智 洋



2.3.3 標準偏差の意味

標準偏差を用いることでデータの中のいったい何がわかるのだろうか。これについては、2つの見方が考えられる。まず第一に、「1セットのデータの中のある1個のデータが持つ意味（特殊性）がわかる」ということ、そして第二に、「複数のデータのセットを比較した場合に出てくる違いがわかる」ということである。

第一の「1セットのデータの中のある1個のデータが持つ意味（特殊性）がわかる」事例を以下に示す¹⁾。数学のテストで自分の結果が75点で平均点が60点だったとする。そのとき、自分はその結果をどれくらい喜んでいいものだろう。このとき知っておく必要があるのは「標準偏差が何点か」ということである。ここでの標準偏差は、テストを受けた人の点数と平均点との差（偏差）を全員分、二乗平均したものである。標準偏差が12点だったとする。そうすると、自分の点数は平均点より「おおそ標準偏差1個分程度高い点数」だとわかる。標準偏差は「各データの平均値からの離れ方を平均化した値」であるから、自分の点数は「平均点から良い方に『普通程度に』離れた値」だということができる。そのような人は大勢いるだろうから、それほど大喜びすることではないといえるかもしれない。

標準偏差がもっと低く、8点だったとする。このとき、自分の点数は「平均的な離れ方である8点」の2倍近くも良い方向に離れていることから、その価値は標準偏差が12点の場合より大きいと評価できるだろう。

以上のように、「1セットのデータ中のある1

個のデータが持つ意味（特殊性）」は、平均値との比較、つまり平均値からの離れ方（偏差）の数値自体だけでは評価できず、標準偏差を基準に相対的に見る必要があるということである。つまり、「そのデータの偏差が標準偏差いくつ分であるか」という視点が重要となってくる。

統計学の世界では、次のようなおおまかな基準が広く了解されているようである。すなわち、1セットのデータの中の「あるデータ」の偏差が、標準偏差 ± 1 個分であれば、そのデータは「月並みなデータ」、そして ± 2 個分から外側のデータである場合は「特殊なデータ」だということである。

続いて、第二の「複数のデータのセットを比較した場合に出てくる違いがわかる」事例を以下に示す¹⁾。例えば、X君は模擬テストを10回受けて平均点が60点で標準偏差が10点、Y君は同じ模擬テストを10回受けて平均点が50点、標準偏差が30点だったとする。このことから何が言えるだろう。

今度は、二人が持つ複数（10回）のデータ（模擬テストの点数）のセットを比較してみる。平均点をみれば、X君の点数の方がY君の点数より優秀である。しかし、これだけでこの2人の本番の受験の結果を見通すことはできない。「標準偏差1個分のバラツキは普通に起こりうる」ということから考えて、X君はおおよそ $60 \pm 10 = 50 \sim 70$ 点の範囲の点数を取る生徒だと見積もることができる。同様に、Y君は $50 \pm 30 = 20 \sim 80$ 点の範囲の点数を取る生徒だと推測できる。つまり、X君は「安定した成績」を取る人で、Y君は「ムラのある成績」を取るタイプだと言える。平均点だけでは評価しきれない側面である。例えば、X君は50点で

入れる学校ならきっと不合格になることはないだろうが、80点取らないと合格しない学校にはなかなか入れないだろう。それに対してY君は、40点で入れる学校にも不合格になる可能性がある反面、80点を要する学校にも合格のチャンスがある、ということもいえるのである。

以上のように、複数のデータセットを比較する際にも、平均値だけでなく標準偏差を考慮すれば別の見方の評価ができることがわかる。

ここまで、データの特徴を表現する「平均値」、「分散」、「標準偏差」という代表的な統計量を説明してきた。ところで、少し抽象的な表現になるが、データというものは世の中で起こった（あるいは将来起こる）さまざまな現象を観測（あるいは予測）することで得られるが、その現象を発現させるシステムの大部分は何かしらの「不確実性」を持っている。そのため、そこから生み出されるデータは判で押したようにぴったり同じ数値になることはほとんどなく、まちまちな値になるのが一般的である。「データがまちまちな数値をとる」ことを専門的には「分布する」という。分布が生じるのは、その数値が決まる背後に何らかの「不確実性」が働いているからであるが、「不確実性」と一口に言っても、そこには固有の「特徴」、「癖」がある。その固有の特徴、癖を専門的には「分布の特性」と呼ぶ。その「分布の特性」を、平均値と標準偏差を用いて「おおよそこの辺を中心にだいたいこれくらいの散らばり具合で存在（分布）する」と表現するわけである。

3. 正規分布

3.1 正規分布の特徴

データの分布について触れたところで、データの分布でもっとも代表的なものと言われる「正規分布」を紹介する。平均値 μ 、分散を σ^2 の正規分布を図3-1に示す。

このグラフは横軸を連続的な確率変数、縦軸を確率密度関数として正規分布をとったものである。確率変数とはいわばそのデータが取りうる数値のことである。それは、サイコロの目(1, 2, 3, 4, 5, 6)のように飛び飛びの値のこともあるが、人の身長や流入下水のBOD値等、切れ目なくどんな

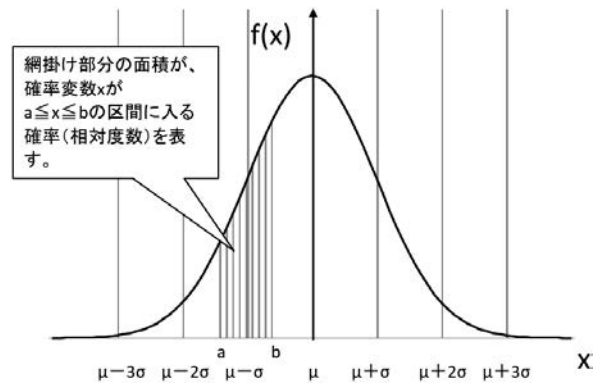


図3-1 平均値 μ 、分散を σ^2 の正規分布

値でも取りうるものは「連続的な確率変数」という。「連続的な確率変数(x とする)に対する確率分布を関数($f(x)$ とする)で表現したもの」が確率密度関数である。この確率分布を表す確率密度関数のグラフでは、横軸とグラフとで囲まれる部分の面積が確率を表す。例えば図4-1で確率変数 x が区間 $a \leq x \leq b$ に入る確率は図4-1に示す「網掛け部分の面積」となる。 x の値から算出される $f(x)$ の値自体が確率を示すのではない。したがって、 x がこのグラフの左端から右端までのすべての範囲の確率(横軸とグラフとで囲まれる部分の面積)は1(100%)となる。

正規分布はこのような確率密度関数の中でも最も有名なものである。その特徴は左右対称な釣り鐘型の分布となっていることである。まず平均値が一番高く、その近辺にもデータが集中して存在し(曲線が盛り上がり)、そこから左右に急激に減少していき平均値 μ から左右 $\pm 2\sigma$ を超えた辺りでは $f(x)$ の値は非常に小さくなる、すべり台のような曲線は正規分布の特徴を示す重要なポイントである。

正規分布は、自然現象や社会現象として観測されるデータに非常に頻繁に現れるものとされている。例えば、人間も含め生物の身長(体長)のデータは正規分布となることが知られている。また、株の収益率のデータも正規分布だと考える研究者は多いということである。

正規分布の姿を示す確率密度関数 $f(x)$ は数式で表すことができる。したがって、先述の区間 $a \leq x \leq b$ をどの範囲で指定しても、そこにデータが入る確率(相対度数)ははっきり決まる(算出

される)。例えば、その中でよく使われるものとして以下のものがある。

- ・平均から標準偏差±1個分の範囲 ($\mu - \sigma \leq x \leq \mu + \sigma$) のデータの相対度数は0.6826 (=70%弱)。
- ・平均から標準偏差±2個分の範囲 ($\mu - 2\sigma \leq x \leq \mu + 2\sigma$) のデータの相対度数は0.9544 (=95%強) (図3-2)。

このことは、前項で説明した「1セットのデータ中のある1個のデータが持つ意味(特殊性)」は、「そのデータの偏差が標準偏差いくつ分であるか」という観点から評価する必要があり、それが±1個分であればそのデータは月並みなデータで、±2個分から外側のデータであれば特殊なデータだ、ということ裏付けている。

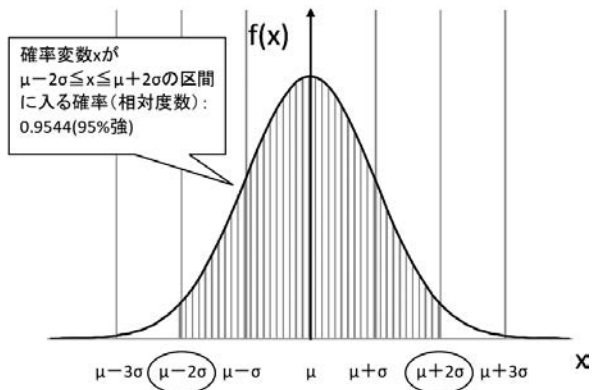


図3-2 平均値 μ 、分散を σ^2 の正規分布

3.2 標準正規分布

標準正規分布とは、平均=0、標準偏差=1の正規分布のことで、正規分布の中で最も基礎になるものである(図3-3)。あるデータ x が平均値 = μ 、標準偏差 = σ の一般正規分布のデータであるとき、

$$z = (x - \mu) \div \sigma \quad \dots \text{(式3-1)}$$

という加工をすると、データ z は標準正規分布のデータに変換される。

<次号へ続く>

【参考文献】

- 1) 完全独習 統計学入門 小島 寛之 2006年9月ダイヤモンド社
- 2) まずはこの一冊から 意味がわかる統計解析 涌井 貞美 2013年2月 ベレ出版

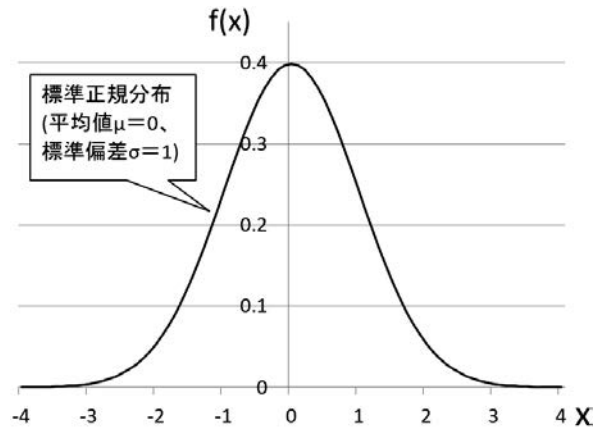


図3-3 標準正規分布

シリーズ・技術調査報告

「統計じかけのオレンジ」

— A Statistics-work Orange —

第3回 推定(1)

一般社団法人 日本下水道施設業協会

技術部長 堅田 智 洋



4. 推定

4.1 統計的推定とは

ある集団に対して何かしらの調査を行うとする。その調査対象となる集団のことを母集団と呼ぶ。その調査は、正確を期すなら対象全部に対して調査を行う「全数調査」が望ましいが、現実には、母集団の規模や費用等の事情から、母集団全部を調査し尽くすことは現実的に不可能なことが多い。そこで、その母集団から一部を選び出し（抽出という）調査することになる。このとき、調査対象として選び出された母集団の一部を「標本」と呼び、標本に対して行う調査を「標本調査」という。

標本調査では、母集団から標本を抽出し、その標本から母数（母集団の分布の特徴を示す、平均値や分散、標準偏差といった数値のこと）に関連する情報を得る。しかし、その調査結果は抽出する標本ごとに異なる。例えば、A市の20歳男子を母集団として、その平均身長を調べるためにその中から10人の男子を抽出したとする。その10の標本の変量（調査項目のこと。ここでは身長）を x_1, x_2, \dots, x_{10} とすると、それらは母集団に従った確率変数（そのデータが取りうる数値）である。この10人の平均身長を \bar{x} とすると、

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_{10}}{10} \quad \dots \text{(式4-1)}$$

となる。このように、「標本から算出された平均値」を標本平均と呼ぶ。重要なことは、この標本平均（ここでは抽出した10人の平均身長）自体も、抽出する標本ごとに値を変える確率変数であるということである。つまり、この10人の標本の平均身長は、抽出する10の標本のセットごとにバラツキが出る。それは、無作為に抽出される男子10人の顔触れがその都度変わるからである。ちなみに、この標本平均のように、母数の推定に利用される統計量をその母数の推定量と呼ぶ。標本平

均とは母平均の推定量である。そして、実際の標本から算出された推定量の値を推定値という。

以上のように、バラツキのある標本調査の結果から、元の母集団に関する真の値（母数）を推定しようというのが「推定（統計的推定）」である。つまり、部分的に明らかになった結果から全体を推測するのである。では、私たちはどのような母集団について、どのような仕組みで統計的推定を行うことになるのだろうか。

例えば、1個のデータ x が現実には観測されたとして、そのことから x が属する母集団について何が言えるだろう。まず、「母平均 μ はこの x の近くだろう」程度の推量は可能であろう（あくまで推量だが）。なぜなら平均値とは分布の中から選ばれた大小いずれにも偏っていない代表的な点だからである。さらに、ここでもし仮に、母集団の標準偏差 σ が何らかの理由でわかっているならどうだろう。先程、1セットのデータの中のあるデータの偏差が標準偏差 ± 1 個分なら、そのデータは「月並みなデータ」で、 ± 2 個分からは外側のデータである場合は「特殊なデータ」と述べた。これは、たいがいのデータは平均値を中心にその前後標準偏差2個以内の範囲におさまることを意味する。つまり、「データ x は母平均 μ から $\sigma \times 2$ 程度以内の離れ方だろう」と考えることができる。逆に言えば「 x から $\sigma \times 2$ 程度以内の隔たりで μ が存在するだろう」、つまり、観測されたデータは、ある程度母平均に近いものとみてよい。このことは、母集団が正規分布の場合には強く支持される。

次に、観測されたデータが1個ではなく、複数個の場合を考える。もちろん、数個のデータでも母集団の分布を再現するほどの情報にはならない。しかし、母平均 μ の推量なら、1個の観測データの時よりもずっと精度を上げることができる。さきほど、複数個の観測データの平均値を母平均

と区別するために標本平均と呼んだが、標本平均については、次の定理が成立する。「1つの母集団から、 n 個のデータを観測しその標本平均 \bar{x} を作る。このとき、 n が大きければ大きいほど、標本平均 \bar{x} は母平均 μ に近い数値をとる可能性が高くなる。」これを「大数の法則」という。「同じ標本平均をとるにしても、より多くのデータでとった方が正確である」という意味のイメージしやすい定理だといえる。

つまり、私たちがこのように標本平均をとるのは、偶然に起きるデータの散らばりをできるだけ打ち消して、実際の値に近い値を得たいからである。この標本平均の発想が統計的推定でも大きな効力をもつ。

さて、統計的推定として最も理想的なものは、もちろん、「母集団の分布がまるでわからないときの推定」だろう。しかし、これは、母集団に関する情報が何もないのだから、さすがに「原理的には無理だ」と直感的に感じさせる。実際には、いくつかのアプローチの方法があるようだが、本稿では取り扱わない。また、正規分布をとらない母集団は、たとえ母集団自体の分布がわかっているとしても、標本平均の分布はそれとは変わってしまうため、推定の対象とするには不都合である。

他方、正規母集団（正規分布している母集団のこと）は、「標本平均を作っても、その分布は正規分布のまま」という都合のいい性質をもっている。そして、私たちが、ある特定の不確实现象に関して何か知りたい場合の多くは、その母集団の母平均や母分散（あるいは標準偏差）を知ればよいことが多い。それは、前述したように、自然現象や社会現象において正規分布が非常に頻りに現れるからである。

したがって、本稿では、統計的推定の対象を正規母集団とする。そして、正規母集団の母平均、母分散をばらつきのある標本調査の結果から推定する方法を下記の4パターンに分類してみていく。

- ・ 正規母集団における母分散 σ^2 がわかっているときの母平均 μ の推定
- ・ 正規母集団における母平均 μ がわかっているときの母分散 σ^2 の推定
- ・ 正規母集団における母平均 μ がわからないときの母分散 σ^2 の推定
- ・ 正規母集団における母分散 σ^2 がわからないときの母平均 μ の推定

4.2 統計的推定のパターン別アプローチ

4.2.1 統計的推定のファーストアプローチ

まず、「区間推定」の重要なポイントは、100%を望まずにピンポイントでなく、例えば「95%の確率で」と幅をもって母集団の平均値や分散値を言い当てようとするところである。本題の「正規母集団とわかっている、母分散もわかっているときの母平均の推定」に入る前に、その区間推定の考え方を、標本数（観測データ）が1個の場合で確認する。たった1個の標本数で区間推定を行うことは実際には現実的ではないが、区間推定のアプローチを学ぶことができる最もシンプルな事例となりうるのでここでみておくことにする。

【例題1】²⁾

お金が入っている箱がたくさん並んでいる。各々の箱の額は不明で、その平均金額を知るために、標本として1箱だけ無作為に抽出して調べたところ、箱の中には500円が入っていた。すべての箱の中の平均金額を区間推定せよ。ただし、箱の中の金額 X は分散 30^2 （標準偏差30）の正規分布に従うとする。

推定するといっても確率現象の話である。「100%こうなる」と断言することは不可能である。そこで、区間推定では、「ここからここまでの幅(区間)に母平均 μ が含まれる確率は95%」というように、推定した値に「当たる確率は何%」と精度を与える。この95%のことを「信頼度」と呼ぶ。得られる推定区間がどれくらい信頼できるかを確率的に表現したものである。区間推定においてよく用いられる信頼度は95%と99%である。以降、本稿では、実際に最もよく使われる信頼度95%でみていくこととする。

推定したい平均金額すなわち母平均を μ とおく。推定量となる金額 X は平均値 μ （未知）、分散 30^2 の正規分布に従うので、抽出した標本ごとに異なる金額で散らばりながら図4-1のように分布する。

前述したとおり、確率分布を表す確率密度関数のグラフでは、横軸とグラフとで囲まれる部分の面積が確率を表す。この図の正規分布でも同じであるので、図4-1で平均値 μ を対称にして左右同じ幅だけ面積（すなわち確率）が信頼度0.95(95%)となるように網掛けをする（図4-2）。

箱の中の金額 X は95%の確率でこの網掛けした区間に現れるはず、と考えるのである。もちろん、この網掛けした範囲の右側と左側に、左右合わせて5

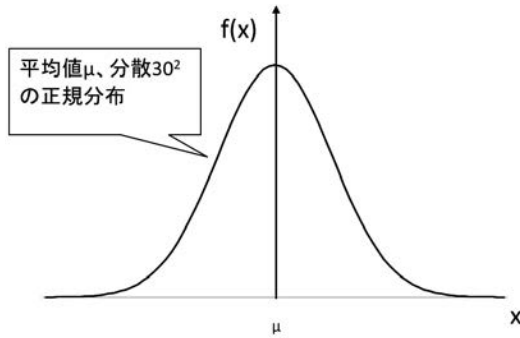


図4-1 平均値 μ 、分散 σ^2 の正規分布

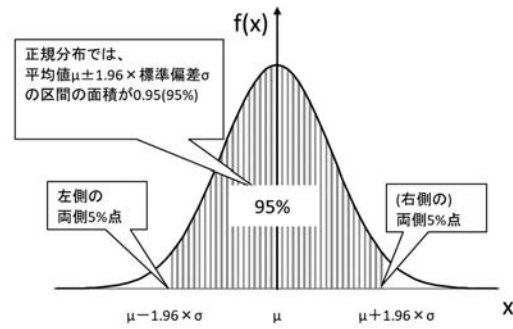


図4-3 平均値 μ 、分散 σ^2 の正規分布における両側5%点

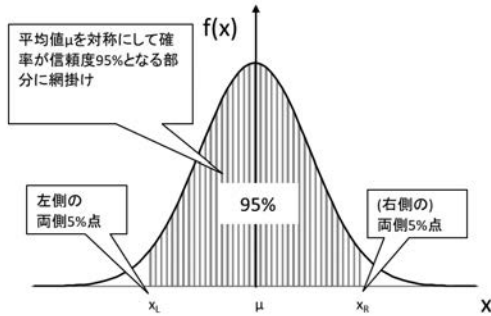


図4-2 平均値 μ 、分散 σ^2 の正規分布 (μ を対称に信頼度95%となる部分に網掛け)

の確率で X が現れる可能性が残っているが、そこも含めた100%の確率で現れる X の区間は $-\infty$ 以上 $+\infty$ 以下になってしまうので、それでは推定する意味がない。そこで、信頼度95%を設定して残り5%の確率で起こり得る事象をカバーすることは諦めることで「有限の区間を推定する」のである。

図4-2の正規分布のように、確率密度変数が左右対称の場合、信頼度95%で網を掛けたときの右端の方の点を「両側5%点」と呼ぶ(x_R)。グラフが左右対称なので、「両側」と言いながら右の方の点のみを指定すれば情報としては足りるのである。ただ、平均値 μ を中心として右側の「両側5%点」 x_R と対称となる x_L を説明したいときには、 x_R を「右側の両側5%点」、 x_L を「左側の両側5%点」と別々の呼称で呼ぶこともある。

平均値 μ 、分散 σ^2 の正規分布の「右側の両側5%点」、「左側の両側5%点」は次のように与えられる(図4-3)。

$$\text{右側の両側5\%点} : \mu + 1.96 \times \sigma \cdots (\text{式4-2})$$

$$\text{左側の両側5\%点} : \mu - 1.96 \times \sigma \cdots (\text{式4-3})$$

よって、網を掛けた X の範囲は次のように表される。

$$\mu - 1.96 \times \sigma \leq X \leq \mu + 1.96 \times \sigma \cdots (\text{式4-4})$$

これは、「正規分布では、 X が95%の確率で平均値 $\pm 1.96 \times$ 標準偏差の区間に出現する」ことを意味する。前節で、正規分布で平均から標準偏差 ± 2 個分の範囲にデータが入る確率は0.9544(95.44%)

であると述べた。そこから95%ぴったりにするために標準偏差 ± 2 個分からほんの少しだけ範囲を狭めて ± 1.96 個分とした、と考えてよい。

この例題では標準偏差 $\sigma = 30$ であるから、

$$\mu - 1.96 \times 30 \leq X \leq \mu + 1.96 \times 30 \cdots (\text{式4-5})$$

とし、金額 X (500円だった)が95%の確率でこの区間に現れるものとして、 μ を逆算するのである。

すなわち、現実を観測した X (500円)がこの範囲に入らないような μ は現実の母集団の母平均としてありえないとして棄却する(採用しない)、というふうに考える。

$$\mu - 1.96 \times 30 \leq 500 \leq \mu + 1.96 \times 30 \cdots (\text{式4-6})$$

から

$$441.2 \leq \mu \leq 558.8 \cdots (\text{式4-7})$$

この範囲に入る母平均 μ は棄却せずに残ることになるので、これが母平均 μ の区間推定の結果、すなわち「母平均 μ の95%信頼区間」となる。

ここで、相対度数が95%になる区間は、(式4-4)の「 $\mu - 1.96 \times \sigma \leq X \leq \mu + 1.96 \times \sigma$ 」以外にもいろいろあるのではないかという疑問が浮かぶ方もいると思われる。その通りである。例えば、(式4-4)の区間を少しずらして「 $\mu - 2.1 \times \sigma \leq X \leq \mu + 1.86 \times \sigma$ 」としても、相対度数は同じ95%になる。しかし、相対度数95%となる区間は、前者の3.92に対して後者は3.96と長くなってしまふ。同じ相対度数で推定する区間は短いほど精度が高いだけ望ましいので、正規分布においては、ヒストグラムが左右対称で、対称軸に近いほど頻度が高いことを考慮すれば、同じ確率で推定される区間を選び出すには、「左右対称の区間」を選択すべきということがわかるだろう。

〈次号に続く〉

【参考文献】

- 1) 完全独習 統計学入門 小島 寛之 2006年9月ダイヤモンド社
- 2) まずはこの一冊から 意味がわかる統計解析 涌井 貞美 2013年2月 ベレ出版