記述統計学 (descriptive statistics)

2008年6月30日

1 はじめに

統計学は、バラツキのあるデータに対して、経験的知識や専門的認識をもとに、数学の手法を用いて、数値 データの特質や規則性の有無を発見あるいは検証します。統計的手法は、統計手法の開発、実験観測の計画、 データの要約や解釈を行う上での理論的な根拠を提供する学問であり、幅広い分野で応用されています。

英語で統計または統計学を statistics といいますが、語源はラテン語で「状態」を意味する statisticum で あり、この言葉がイタリア語で「国家」を意味するようになり、国家の人力、財力等といった国勢データを比 較検討する学問を意味するようになったといわれています。。現在では、経済学、自然科学、社会科学、医学 (疫学、EBM)、薬学、心理学、言語学など広い分野で必須の学問となっていることは論をまたない。また統計 学は哲学の一分科である科学哲学においても重要なひとつのトピックスになっている。理由は統計学が科学的 な研究においての方法論として基礎的な部分を構成していながら、確率という極めて捉えがたい概念を基盤と していることもあり、不確実性の知覚についてその意味やあり方が帰納の正当性の問題などと絡めて議論され ている。

記述統計とは、収集したデータから要約した数値、すなわち統計量とよばれる平均、分散などの数値を計算 して、データの状況や規則性を明らかにすることで、データの示す傾向や性質を知ることです。とくにデータ については、これに関する経験や知識が分析を行うためには重要なかぎとなる。新たな発見、知覚認識に結び ついていく。

推測統計という分野では、抽出されたデータからその根源となっている諸性質を確率論的に推測する分野で あり、経験から発見される知覚ではなく、数理的な仮説にもとづく検証や推測を目的とする。

2 統計調查

新聞やテレビにも多くの統計調査が発表されている。いくつかの代表的な統計調査を挙げてみる。

(i) 政府による官庁統計:人口動態(出生率、死亡率、市町村別人口動態)、消費者物価指数、GDP 成長率

- (ii) マスコミによる調査: 世論調査、政党支持率、番組視聴率調査
- (iii) 民間会社による市場調査: POS データ、商品動向、需要予測、経済予測

これらの調査ではいくつかの項目を調べる。

調査における主要な項目:調査の企画、目的と対象、経費、実施手順、調査項目、設計、抽出、2段階層別 抽出法調査結果の処理、集計、データ処理、結果の分析、公表

3 統計データの種類、分類

尺度水準(しゃくどすいじゅん)とは、調査対象に割り振った変数、その測定、あるいはそれにより得られ たデータを、それらが表現する情報の性質に基づき数学・統計学的に分類する基準である。データ(あるいは 変数、測定)の尺度あるいは単位の構造から、ふつう次のような種類(水準)に分類される。この尺度水準に よって、統計に用いるべき基本の統計量や統計検定法が異なることに注意する。

3.1 質的データ

質的データは、カテゴリデータともよばれ、観測データとしては選択あるいは観測した結果が項目(カテゴ リー)である。

● 名義尺度:単なる番号であり、順番の意味をもたない。電話番号、背番号など。

この水準では数字を単なる名前として対象に割り振る。2つの対象に同じ数字がついていればそれらは同じ カテゴリに属する。変数値間の比較は等しいか異なるかでしか行えない。順序もないし加減などの演算もでき ない。例としては電話番号、背番号、バスの系統番号など。中心的傾向の指標として使えるのは最頻値のみで ある。統計的バラツキは変動比や情報エントロピーで評価できるが、標準偏差などの概念はありえない。名義 尺度でのみ測定されるデータはカテゴリデータとも呼ばれる。

* なおカテゴリデータを、ある性質が「あるかないか」という表現に直し、さらにこれを「1か0か」で表現 したものをダミー変数という。ダミー変数またはそれから算出されるスコア(点数)を、順序尺度以上の水準 に準じて扱う方法もよく用いられる。

● 順序尺度:順序が意味を持つ番号。階級や階層など。

この水準では対象に割り振られた数字は測定する性質の順序を表す。数字は等しいかどうかに加え、順序 (大きいか小さいか)による比較ができる。しかし加減などの演算には意味がない。物理学的な例にはモース 硬度がある。その他の例にはレースの着順などがあるが、これでは到着時間の差は記録できない。心理学や社 会科学の測定のほとんどは順序尺度で行われる。例えば社会的態度(保守的か進歩的かなど)や階級は順序水 準で測定されるものである。また客の嗜好(アイスクリームのバニラ味とチョコレート味とどちらが好きか) のデータもこれで表現できる。順序尺度の中心的傾向は最頻値や中央値で表されるが、中央値の方が多くの情 報を与える。順序尺度で測定されるデータは順序(または順位)データと呼ばれる。

* 以上の名義尺度および順序尺度で表されるデータを合わせて質的データともいう。また各カテゴリに属す 対象の個数という形のデータにまとめると数量データと呼ばれ、これは分割表で表示できる。これらに対して 用いられる統計検定法はノンパラメトリックなものに限られる。

3.2 量的データ

いわゆる観測、調査の結果えられたものが数値データであるものをさす。尺度構造として、間隔尺度と比率 尺度に分ける。また変数値が実数と整数のそれぞれの場合について、実数型と離散型データとよばれる。

● 間隔尺度:順序に加え間隔にも意味がある(単位がある)が、ゼロには絶対的な意味はない。摂氏・華氏 温度、知能指数など。

対象に割り振られる数字は順序水準の性質を全て満たし、さらに差が等しいということは間隔が等しいとい うことを意味する。つまり測定値のペアの間の差を比較しても意味がある。加減の演算にも意味があるが、尺 度上のゼロ点は任意で負の値も使える。例にはカレンダーの日付がある。値の間の比には意味がなく、直接の 乗除の演算は行えない。とはいえ差の比には意味がある。中心傾向は最頻値、中央値あるいは算術平均で表さ れ、算術平均が最も多くの情報を与える。間隔尺度で測定されるデータは間隔データと呼ばれる。摂氏または 華氏で測る温度も間隔尺度である。社会・人文科学分野で普通用いられる唯一の間隔尺度は知能指数(IQ)で ある。

●比率尺度:ゼロを基準とする絶対的尺度で、間隔だけでなく比率にも意味がある。絶対温度、金額など。 対象に割り振られた数字は間隔尺度の性質を全て満たし、さらにその中のペアの比にも、乗除の演算にも意 味がある。比率水準のゼロ点は絶対的である。ほとんどの物理学的量、つまり質量、長さやエネルギーは比率 水準である。また温度も絶対温度で測れば比率尺度である。比率尺度で測定される変数の中心的傾向は最頻 値、中央値、算術平均あるいは幾何平均で表されるが、間隔尺度と同じく算術平均が最も多くの情報を与える。 比率尺度で測定されるデータは比率データと呼ばれる。比率尺度で表される社会的変数には年齢、ある場所で の居住期間、収入などといったものがある。

* 正しい意味で単位を有するのは間隔尺度と比率尺度のみであり、従ってこれらは真の尺度とも呼ばれる。 これらのデータを合わせて量的データ(質的データに対して) 数値データ(数量データに対して)ともいう。

スタンレー・スティーヴンズ (Stanley Smith Stevens)により 1946 年の論文「測定尺度の理論につい て」"On the theory of scales of measurement"で提案された分類がよく用いられる。変数に対して可能な数 字の演算は、変数を測定した尺度水準に依存し、その結果、特に統計学で用いるべき要約統計量および検定法 も変数の尺度水準に依存する。スティーヴンズは低い方から順に以上の4つの尺度水準を提案しており、高い 水準はより低い水準の性質を含む形になっている。また高い水準でのデータを低い水準に変換して扱うことが できる。参考:Wikipedia, Excel help

4 度数分布表とヒストグラム

度数分布(どすうぶんぷ、Frequency Distribution)とは、統計において標本として得られたある変量の値の リストである。一般に量の大小の順で並べ、各数値が現われた個数を表示する表(度数分布表)で示されます。

例。ある30人のクラスで、体重の測定値は量的データとして得られます。

例。例えば、100人がある文章に同意するかを5段階のいわゆるリッカート尺度で回答したとします。これ は項目によって分類され、質的データでの順序尺度の単位構造をもちます。このとき、1は強く同意すること を示し、5は全く同意しないことを示す。その回答群を度数分布で表すと次のようになります:

階級	同意の度合	回答数	
1	強く同意する	25	
2	ある程度同意する	35	
3	どちらとも言えない	20	
4	ある程度同意できない	15	
5	全く同意できない	5	

ヒストグラム(柱状グラフ)

この単純な表には2つの弱点がある。変量が連続的な 値をとりうる場合や非常に範囲が広い場合、度数分布 表の作成は難しくなる。

連続的な変量に対しては、観測データをいくつかの階 級に分類し、その中心となる階級値ごとに度数を修正 する。つぎに述べるヒストグラム(Histogram)(柱状 グラフとのよばれる)と関連されて表現する。

階級によって区分けされた観測データをグラフによって表現します。注意すべきことはこの階級の幅に意味 をもつことです。それに対して棒グラフで表す場合ではデータが離散的な場合やカテゴリーで与えられている 場合であって、階級に区分けするときは、項目の数(x軸に対応する)が多数であってある程度まとめたほう がよい場合です。

平均と中央値が異なる場合、度数分布に歪み(ひずみ)があるといいます。正規分布は平均を中心として対称な形状をして歪みがないといいますが、これとは異なり、L 字型と逆 L 字型の場合が相当します。L 字型とはデータの値が(平均値の右側)大きいほうに低く長くあることから、3次モーメントは正の値となり、この計算から歪度は正値となります。逆に大きい値により多くのデータがかたまり、この付近に度数も大きければ、平均もこの大きいほうになり、平均の周りの3次モーメント値は(平均値の左側)小さい値が多数になることから、逆 L 字型の分布では歪度が負の値になります。

度数分布の尖度(せんど)とは、平均値への集中の度合であり、ヒストグラムで表した場合のグラフの尖り (とがり)具合です。正規分布以上に尖っている場合を「急尖的;leptokurtic」と称し、逆の場合を「緩尖的; platykurtic」とよばれます。データのばらつき具合を正規分布と比較するもので、数値3が正規分布であり、 平べったいならば、平均の周りの4次モーメントを計算することから、値が大きくなり、負の値になる場合は 尖っていて、平均の周りにデータが固まっていてモーメントの値が小さくなるからです。

ヒストグラム(度数分布図、柱状グラフ、Histogram)とは、縦軸に度数、横軸に階級をとった統計グラフ の一種で、データの分布状況を視覚的に認識するために主に統計学や数学、画像処理等で用いられる。

累積度数分布表

級中央値に対する度数の対応表は度数分布表ですが、度数を累積していったものは、単調に増えていき、最 後の値はデータの総数になります。これは確率分布では分布関数に相当します。中央値、四分位数、十分位数、 百分位数などの計算はこの累積度数分布表から、それぞれ、50%づつ2つに分ける、25%づつ4つに分け る、10%づつ10個に分ける、1%ごとの100個のデータに大きさの順にしたものとなります。

• 相対度数分布表

上で述べた累積相対度数に対して、データの総数で割り、百分率、パーセント値に直したもの。このように 基準化することで、データ数が多数であっても表示され、2つ以上の比較も可能となります。

統計図表

いわゆる統計グラフは、統計図表(とうけいずひょう)ともよばれ、複数の統計データの整理、視覚化、分 析、解析等に用いられるグラフおよび表の総称を意味します。ここで、グラフとは「図形を用いて視覚的に、 複数の数量・標本資料の関係などを特徴付けた物」のことを指します。この意味においてのグラフはしばし 「統計グラフ」と呼ばれます。統計図表は、統計データの整理、分析、検定などの過程で用いられる。統計図表 を駆使することで、「調査活動によって得られた数量(統計データ)の特徴」(増減の傾向の型,集団の構成な ど)や、統計データ同士の関係(相関関係など)を視覚的に理解することが出来ます。

統計グラフの種類

統計グラフには、様々な種類がありますが、以下に典型的な統計グラフを示します。

- 棒グラフ 棒グラフは、資料を質的に(意味的に複数の項目に)分類したときに、各項目間の大きさを比較する ために用いる。項目を横軸、各項目の大きさを縦軸に表現する(横軸、縦軸は逆でも良い)。棒で表すこ とで、各項目の大きさや、大きい値(小さい値)を持つ項目、各項目間の関係などが把握しやすくなる。
- 柱状グラフ(ヒストグラム) 柱状グラフ(ヒストグラム)は、棒グラフの一種で、資料を量的に(大きさを 複数の階級に区分し、各要素がどの階級に属するかという指標で)分類した時に、各階級の散らばりの 様子を見るために用いる。柱状で表すことで、集団の偏りや各階級間の散らばりの様子が把握しやすく なる。品質管理などにおいて、度数分布表から度数分布を図示するときによく用いられる。度数が増え るにしたがって、グラフの形状は柱状から曲線へと近づいてゆく。この曲線を度数分布曲線という。
- 円グラフ 円グラフは、資料を特定の項目に分類した時、その一項目での割合を比較する時によく用いられる。 円で全体を表すことで、ある項目内・分野内での割合の大小が直感的に把握しやすく、プレゼンテー ションなどでよく利用される。又、円グラフでは、全体の数値を360として表現することも少なくない。 他方で、厳密な比較には向かないため、専門分野ではむしろ使用されない。
- 折れ線グラフ 時刻変化を表すためには、時刻に対応した数値を直線で結んでいく。また2つ以上の系列グラ フを比較にもよく用いられる。
- 絵グラフ 視覚的な強調や興味あるアピールを行うために様々な表現方法で表示する。一般に表示を工夫した もの、分類できないようなものなどが含まれる。
- レーダー・チャート 円形図形をもちいた各項目の数値を中心から半径への距離に基準化して表す。多変量 データをまとめて表せる。たとえば、各人のテストでの科目ごとの得点結果をまとめて表現したり、

グラフ・ウィザートをもちいて、度数分布表のつくり方を説明します。

- 1. 得られたデータの配列から、最大 (max) と最小 (min) を求め、まず範囲 (range) を計算します。
- 2. 階級 (class) の数をおおよそ 10 から 20 ぐらいになるよう設定します。もしうまく行かない場合は数を

代えてみます。Starjes の公式 $k = 1 + \log_2 n = 1 + 3.3 \log_1 0n$ という目安もあります。

- 3. 階級幅 (class width) が範囲を階級の数で割ること *R/k* で計算できます。階級の度数によってはいくつ かの階級を合併してまとめたほうがよいこともあります。
- 4. これから階級、階級値、度数、累積度数、相対度数、累積相対度数をまとめて表計算ソフトで計算し ます。

	А	В	С	D	Е	F	G
1	階級(下限)	階級 (上限)	階級値	度数	累積度数	相対度数	累積相対度数
2	$a_0 \sim$	a_1	x_1	f_1	$F_1 = f_1$	$p_1 = f_1/n$	$P_1 = p_1$
3	$a_1 \sim$	a_2	x_2	f_2	$F_2 = f_1 + f_2$	$p_2 = f_2/n$	$P_2 = p_1 + p_2$
4	$a_2 \sim$	a_3	x_3	f_3	$F_3 = f_1 + f_2 + f_3$	$p_3 = f_3/n$	$P_3 = p_1 + p_2 + p_3$
5			•••				
6	$a_{k-1} \sim$	a_k	x_k	f_k	$F_k = n$	$p_k = f_k/n$	$P_k = 1$
7	計			n		1	
$x_i = (a_{i-1} + a_i)/2, (i = 1, 2, \dots, k)$ $F_k = f_1 + f_2 + \dots + f_k = n, P_k = p_1 + p_2 + \dots + p_k = 1$							

この表の作成には数式メニューに組み込み関数をいれていきます。=max(データ範囲)、=min(データ範囲) から範囲をあらかじめ計算しておきます。階級の幅(上限と下限)や個数を決めます。度数の枠にいれる命 令は、セル D2 からセル D6 をアクティブに選んでから、=frequency, データ配列、区間配列(セル A2 から セル A6)、入力は数式配列ですから、CTRL+SHIFT+ENTER とします。階級の入力は「編集」 「連続 データの作成」、また累積度数はセル E2 に「=D2」さらにセル E3 に「=D2+D3」として、E3 を E4 から E6 までコピーペーストをしていきます。相対参照をしながら自動計算してくれます。相対度数はセル F2 には 「=D2/D\$7」といれて、分母には絶対参照するセルの値、すなわち合計で割れば求まります。この度数分布表 から、グラフ・ウィザードで棒グラフを選び、ヒストグラムをつくります。系列で×軸に表すラベルには階級 値を指定し、とくにオプション・コマンドで棒グラフの間隔幅をゼロに指定します。あるいはこの frequency 命令ではなく、分析ツールのヒストグラムでも作ることができます。

統計グラフにはどのような選択の目安を下記に示す。

- 1.2種類の系列(2変量データの全体を俯瞰する)からなるデータの相関 散布図
- 2.1種類の系列(1変量データで時間とともに変化する)からなるデータの時間的推移(時間との相関)
 2つ以上のグラフ間における比較 折れ線グラフ
- 3. 値の大きさの比較、横軸には項目カテゴリー 棒グラフ
- 4. 集団における内訳や構成比を見る 円グラフ

5 データの代表値

データの集まりから、この分布を特徴づける基本統計量として平均、標準誤差、中央値(メジアン) 最頻値 (モード) 標準偏差、分散、尖度、歪度、範囲、最小、最大、合計、標本数が表示されます。

- 1. 最初の準備として Excel の「ツール」 「アドイン」 「分析ツール」と「分析ツール-VBA 関数」に チェックがあることを確認
- 2.「ツール」 「分析ツール」 「基本統計量」 「OK」
- 3. 基本統計量のメニューでは入力範囲に「データの範囲」を指定し、データの先頭行からすべてデータ値のときには「先頭行をラベルとして使用」のチェックマークをはずす。出力オプションは「新規またはつぎのワークシート」になっていますが、もし出力先を同じシートにする場合には、セルを指定します。 基本統計量を出力するためには必ず「統計情報」にチェックを入れます。最後には「OK」とします。

出力された基本統計量の説明をします。

$$X_1, X_2, \cdots, X_n$$

とします。これを大きさの順に並び替え(ソート,sort)したものを順位統計量 (Order Statistics) といい、 $X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(n)}$ と表す。これを階級別に整理されたとき (k は固定された階級の数) には

階級値 (変数)	x_1	x_2	 x_k	計
度数	f_1	f_2	 f_k	n
相対度数	p_1	p_2	 p_k	1

の形でデータが与えられたとします。大文字と小文字でこの2通りを区別していますが、 $f_i = 1, i = 1, 2, \cdots, n$ とすれば、後者は前者の場合に帰着されます。これはそれぞれのデータがすべて異なった値をとっていることになります。

$$X_1 + X_2 + \dots + X_n = x_1 * f_1 + x_2 * f_2 + \dots + x_k * f_k,$$

$$f_1 + f_2 + \dots + f_k = n$$

$$X_1^2 + X_2^2 + \dots + X_n^2 = x_1^2 * f_1 + x_2^2 * f_2 + \dots + x_k^2 * f_k,$$

$$\frac{1}{n}(X_1 + X_2 + \dots + X_n) = x_1 * p_1 + x_2 * p_2 + \dots + x_k * p_k,$$

$$\frac{1}{n}(X_1^2 + X_2^2 + \dots + X_n^2) = x_1^2 * p_1 + x_2^2 * p_2 + \dots + x_k^2 * p_k$$

平均 算術平均を計算します。 $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i = \frac{1}{n} \sum_{j=1}^{k} x_j f_j = \sum_{j=1}^{k} x_j p_j$ AVERAGE 関数の値と同じ。 標準誤差 SE $SE = \frac{s}{\sqrt{n}} = \sqrt{\frac{\sum_i (X_i - \overline{X})^2}{n(n-1)}} = \sqrt{\frac{\sum_i X_i^2 - n\overline{X}^2}{n(n-1)}} = \sqrt{\frac{\sum_j x_j^2 f_j - n\overline{X}^2}{n(n-1)}} = \frac{\sqrt{\sum_j x_j^2 p_j - \overline{X}^2}}{\sqrt{n-1}}$ 中央値 (メジアン) Me データを大きさの順に並べたとき、中央に位置する値。 $Me = X_{((n/2)+1)}$ (偶数) 、 $= X_{((n+1)/2)}$ (奇数) データ数が奇数・偶数によって真ん中の値がかわる。数式メニューでは = median(データ範囲)、あるいは四分位数の 2 番目の値 = quartile(データ範囲, 2)最頻値 (モード) Mo データの中で、最も頻度が高く (最大度数) 現れた値。 MODE 関数と同じ。もし頻度

の最大が2つ以上で同じ値を取るとき(つまり最大値がないとき)は、"#N/A"(Not Available)と表示される。度数の値 $\max_i f_i$ が最大となるような変数の値。データの中で、最も頻度が高く現れた値。 MODE 関数と同じ。単峰ならばひとつに定まるが、もしひとつに定まらないような双峰形では存在しないとする。データがすべて、異なる値を取るときは、"# N/A"(Not Available)と表示される。

標準偏差 SD $SD = s = \sqrt{\frac{\sum_i (X_i - \overline{X})^2}{(n-1)}}$ 標本分散の平方根を取ったもの。STDEV 関数と同じ。VAR 関数の平方根を取ったもの。

分散(標本不偏分散)s² VAR 関数と一致。

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{i} - \overline{X})^{2} = \frac{1}{\binom{n}{2}} \sum_{i < j} (X_{i} - X_{j})^{2} = \frac{n}{n-1} \left(\sum_{j=1}^{k} x_{j}^{2} p_{j} - \overline{X}^{2} \right)$$

VARP 関数は、分母が(n-1)でなく、nを用いている。

$$v^{2} = \frac{1}{n} \sum_{i=1}^{n} (X_{i} - \overline{X})^{2} = \frac{1}{n} \sum_{j=1}^{k} x_{j}^{2} f_{j} - \overline{X}^{2} = \sum_{j=1}^{k} x_{j}^{2} p_{j} - \overline{X}^{2}$$

尖度(せんど)KW KURT 関数で与えられる量。 $\frac{1}{n}\sum_{i}\left(\frac{X_{i}-\overline{X}}{s}\right)^{4}$ -3とがり具合を測る尺度となります。 歪度(わいど)SK SKEW 関数で与えられる量。 $\frac{1}{n}\sum_{i}\left(\frac{X_{i}-\overline{X}}{s}\right)^{3}$ 左右対称性を測る尺度となります。た とえば、L 字型の形状の分布では正の値をとり、逆 L 字型では負の値をとります。分布が L 字型であれ ば、代表値の大きさは モード < 中位数 < 平均 の順になり、もし逆 L 字型では、順序が、 平均 < 中 位数 < モード の順になります。

範囲 R データの最大値から、最小値を引いたもの。 $R = \max(データ範囲) - \min(データ範囲)$

- 最大 データの最大値。MAX 関数。
- 最小 データの最小値。MIN 関数。
- 合計 データの合計値。SUM 関数。
- 標本数 データ数 n。COUNTA 関数で与えられる値。
- 信頼区間(95.0%) 信頼係数 95% に対する信頼区間の幅の 1/2。上の平均プラスこの値と平均マイナスこの 値で出来る区間が信頼係数 95% の信頼区間を構成する。

上記は基本統計量による出力ですが、これ以外にも重要な統計量があります。

トリム平均 TM(trimmed mean) データに異種的な要素が加わっていると、平均値は大きく影響を受ける。 これを除去して残りのデータについて平均値をとることがある。切り落し平均ともよばれる。データ全 体の上限と下限から一定の割合のデータを切り落とし、残りの項の平均値を返します。TRIMMEAN 関数は、極端な観察データを分析対象から排除する場合に利用します。 書式 TRIMMEAN(配列,割合):(1) 配列 対象となるデータを含む配列またはセル範囲を指定しま

す。(2) 割合 平均値の計算から排除するデータの割合を小数で指定します。たとえば、全体で 20 個 のデータを含む対象に対して割合に 0.2 を指定した場合、20 × 0.2 = 4 となり上限から 2 個、下限か ら 2 個の合計 4 個のデータが排除されることになります。

- 平均偏差 MD (Mean Deviation) = $\frac{1}{n} \sum_{i} |X_i \overline{X}|$ 各データと平均値との差に対する絶対値の和であるが、 このように絶対値を考えることはデータと平均値との距離であり、ばらつきの尺度としてはごく自然な ものである。
- Z スコアー (偏差値) $Z_i = \frac{X_i \overline{X}}{s_x}$ 一次変換により、シフト (平均値の移動) とスケール (分散の拡大縮小) により平均値が 0 となり、分散あるいは標準偏差を 1 に変換する。さらに平均 50、標準偏差 10(分散 100) にしたもの、 $10Z_i + 50$ を変量 X の偏差値という。知能指数は平均 100、標準偏差 15 にしたもの。
- 四分位数、四分位偏差 Q₁, Q₂, Q₃ とは、分布全体を 25% ずつの 4 つに分けて、最小値から 25% になるもの Q₁ を第一四分位数、第 2 四分位数は中央値と一致する 50% のところ、すなわちちょうど半分のところ である。第 3 四分位数、75% のところは第 3 四分位数とよばれる。

数式入力では=QUARTILE(配列, 戻り値)となる。ただし 戻り値 QUARTILE 関数の戻り値

0	最小値
1	第 1 四分位数 (25%)
2	第 2 四分位数 = 中位数 (50%)
3	第 3 四分位数 (75%)
4	最大値

四分位偏差とは $Q_3 - Q_1$ であり、箱ひげ図(ボックスチャート)にも表示される。これ以外にもデー タ数が大きければ、十分位数 (decitile) や百分位数 (percentile) などが用いられる。箱ひげ図(ボックス チャート)では最小値、第1四分位数、中央値(第2四分位数)最大値の4つの値ももちいて表示する。 変動係数 CV(Coefficient of Variation)= $\frac{s_x}{\overline{v}}$

標本共分散
$$Cov(x,y) = \frac{1}{n-1} \sum_{i} (X_i - \overline{X})(Y_i - \overline{Y}) = covar(x 範囲, y 範囲)$$

標本相関係数 $\rho = \frac{Cov(x,y)}{s_x s_y} = Cov\left(\frac{x}{s_x}, \frac{y}{s_y}\right) \exists z = \tau \frac{x}{s_x}, \frac{y}{s_y} \exists \frac{X_i - \overline{X}}{s_x}, \frac{Y_i - \overline{Y}}{s_y}, i = 1, 2, \cdots, n \ge 0,$
 $z = n \varepsilon = \frac{1}{2} \varepsilon$

スピアマンの順位相関係数 変量の値が順位で与えられているとき、 $r_s=1-rac{6\sum_i (x_i-y_i)^2}{n(n^2-1)}=1-$

 $\sum_{i} (x_i - y_i)^2 / {\binom{n-1}{3}}$ ここでデータ x_i, y_i は 1,2,...,n のうちのいづれか一つとなっている。もし 同順位のある場合では、補正因子を入れて、それらの平均順位として調整する。同順位の補正因子 =[COUNT(範囲) + 1 - RANK(数値, 範囲, 0) - RANK(数値, 範囲, 1)]/2 ただし RANK(数値, 範囲, 順序) で、順序:数値の順位を決めるため、範囲内の数値を並べ替える方法を指定します。順序に 0 を 指定するか、順序を省略すると、範囲内の数値が ...3、2、1 のように降順に並べ替えられます。 順序に 0 以外の数値を指定すると、範囲内の数値が 1、2、3、... のように昇順で並べ替えられます。

6 相関係数

2 変量の度数分布表、2 変量のヒストグラムにおいて、相関図(そうかんず)または散布図(さんぷず)と は、縦軸、横軸に2項目の量や大きさ等を対応させ、データを点でプロットしたものである。各データは2項 目の量や大きさ等を持ったものである。

散布図の例:グラフ・ウィザードで散布図を選択する。

散布図には、2項目の分布、相関関係を把握できる特長がある。データ群が右上がりに分布する傾向で あれば正の相関があり、右下がりに分布する傾向であれば負の相関がある。相関係数が0であれば無相関と なる。

共分散: covar(配列 A, 配列 B)

ー連の個別の対象物に対して測定される N 個の異なる測定変数がある場合、相関分析ツールと共分散分析 ツールは同じ設定で使うことができます。相関分析ツールと共分散分析ツールは共に、測定変数の各組み合わ せ間のそれぞれ相関係数または共分散を示すマトリクスが、出力テーブルとして得られます。相関係数が -1 から +1 までの範囲に収まるのに対し、対応する共分散はこの範囲に収まらない点が異なります。相関係数 と共分散は共に、2 つの変数が一緒に変化する範囲で測定されます。共分散分析ツールは測定変数のそれぞれ の組み合わせについて COVAR ワークシート関数の値を計算します。たとえば N=2 の 2 つの測定変数のみ の場合は、共分散分析ツールではなく COVAR 関数を直接使用する方法が適しています。共分散分析ツール の出力テーブルで対角線上の i 行と i 列の値は、それ自身の i 番目の測定変数の共分散を表します。これは、 VARP ワークシート関数で計算されるその変数に対する母集団の分散の値と同じです。共分散分析ツールを 使うと、測定変数の組み合わせそれぞれについて 2 つの測定変数が一緒に変化する傾向があるかどうかを調 べることができます。一方の変数の大きな値がもう一方の変数の大きな値と関連する傾向があるか (正の共分 散)、一方の変数の小さな値がもう一方の変数の大きな値向があるか (負の共分散)、両方の変数 の値が関連しない傾向があるか (0 に近い共分散) などを調べることができます。

相関係数: CORREL(配列 A, 配列 B)

N 個の対象物それぞれに対して各変数の測定を行う場合、CORREL ワークシート関数と PEARSON ワー クシート関数は共に 2 つの測定変数間の相関係数を計算します。いずれかの対象物に対する観察が行われない と、分析時にその対象物が無視されます。相関分析ツールは、N 個の対象物それぞれに対して 3 つ以上の測定 変数がある場合に特に役立ちます。この分析を行うと、測定変数の可能な組み合わせそれぞれに対して適用さ れた CORREL (または PEARSON) 関数の値を示した相関マトリクスが、出力テーブルとして得られます。 共分散と同じように、相関係数は 2 つの測定変数が一緒に変化する範囲で測定します。共分散とは異なり、相 関係数は 2 つの測定変数を表現する単位とは関係なくその値の基準が決められます。たとえば、2 つの測定変 数が重量と高さの場合、重量がポンドからキログラムに変更されても相関係数の値は変わりません。相関係数 のすべての値は、-1 から +1 までの範囲に収まる必要があります。相関分析ツールを使うと、測定変数の組み 合わせそれぞれについて 2 つの測定変数が一緒に変化する傾向があるかどうかを調べることができます。一方 の変数の大きな値がもう一方の変数の大きな値と関連する傾向があるか (正の相関)、一方の変数の小さな値が もう一方の変数の大きな値と関連する傾向があるか (負の相関)、両方の変数の値が関連しない傾向があるか (0 に近い相関) などを調べることができます。

7 回帰分析

回帰分析ツールは、線形回帰分析を行います。回帰分析では、R-2 乗値を使って、観測値のデータが最適な 直線に当てはめられます。このツールを使って、複数の独立変数が1つの従属変数に与える影響を分析するこ とができます。たとえば、スポーツ選手の年齢、身長、体重などの要素が成績に与える影響を分析できます。 成績データに基づいて、これらの要素それぞれが成績に影響した比率を割り当てたり、回帰分析の結果を使っ て、ほかのスポーツ選手の成績を予測することもできます。回帰分析ツールは LINEST ワークシート関数を 使用します。

8 インターネットによる統計データ資料

総務省統計局

http://www.stat.go.jp/data/index.htm

千葉県の統計情報

http://www.pref.chiba.jp/outline/statistics/index-j.html