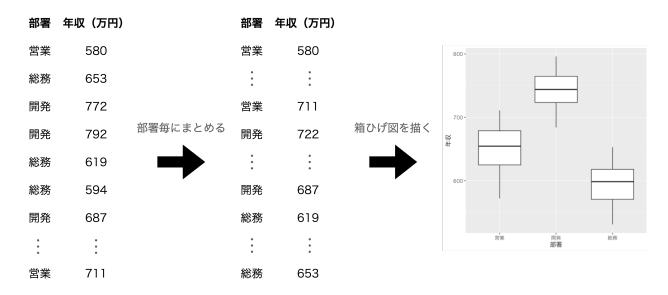
データを読む

データの整理

データから有益な情報を引き出すためには、まず最初にデータを体系的な方法に則り整理するとが必要となります。整理をしなければデータはただの数字や文字などの羅列であり、そこから傾向や特徴を読み取ることが困難な為です。

例えば、ある企業の社員100名分の部署と年収がペアになったデータのリストがあったとします。部署と年収のペアはランダムに並んでいるようで、眺めているだけではこのデータからはなにもわかりません。しかし、部署毎に年収をまとめ、箱ひげ図(後述)を描いてみると、部署によって年収が異なっており、開発>営業>総務の順に高額となっているらしいことが見えてきます



このようにデータを整理することが、データから情報を読み取る足がかりとなるのです.

データの種類

データの整理の仕方はデータの種類によって異なってきます。この節ではデータの種類について見ていきます。

データは質的データと量的データに大別されます。質的データとは分類を表すデータで、算術計算することはできないデータのことです。一方、量的データとは数値で表されたデータで、算術計算することができるデータのことです。

質的データの例としては学籍番号や氏名,順位などが挙げられます.数字で表されたデータであっても計算結果に意味がないデータであれば、それは質的データとなります。Aさん、Bさん、Cさんの学籍番号がそれぞれ「10010」、「10020」、「20030」であるとき、20030(Cさんの学籍番号)= 10010(Aさんの学籍番号)+10020(Bさんの学籍番号)となりますが、「CさんはAさんとBさんを足し合わせたような人物」とはなりません。計算結果には意味がありません。よって、学籍番号は数字で表されていますが質的データととなります。

量的データの例としては気温や西暦、身長などが挙げられます。量的データは数字により表現され、計算結果には意味が伴います。ある日の札幌市、那覇市の気温がそれぞれ「15°C」、「30°C」であるとき、30°C - 15°C = 15°Cと引き算することで札幌市と那覇市の気温差は15°Cであることがわかります。計算により意味のある結果を得ることができます。よって、気温は量的データとなります。

量的データは質的データに変換することができます。量的データである気温に 15° C未満は「寒い」, 15° C以上 25° C未満は「快適」, 25° C以上は「暑い」と温熱感覚を対応付ければ質的データとして扱うことができるようになります。しかし,温熱感覚が「快適」となっていても,この快適の気温を一意に決めることができないので質的データを量的データに変換することはできません。

気温(量的データ)	温熱感覚	(質的データ)
15°C		快適
30°C	対応付ける	暑い
31°C	_	暑い
18°C		快適
•		•
12°C		寒い

質的データおよび量的データに大別されたデータは表現する性質を基準にして更に分類することができます。この基準のことを尺度と呼び、4つの水準「名義尺度」、「順序尺度」、「間隔尺度」、「比例尺度」に分類されます 1 質的データは名義尺度のデータと順序尺度のデータに分けられ、量的データは間隔尺度と比例尺度に分けられます。

名義尺度とはデータの区別(だけ)が可能な尺度です。名義尺度のデータには便宜的に数字を割り振ることができますが、割り振った数字の算術計算には意味がありません。また、大小関係にも意味がありません。しかし、区別はできなければなりませんので等しいか等しくないかには意味があります。

名義尺度のデータの例としては学籍番号を挙げることができます。学籍番号は学生一人一人に数字を割り振った記号として捉えることができます。この学籍番号が等しければ同じ学生を、等しくなければ異なる学生を示していることになります。しかし、先に見たように学籍番号の算術計算には意味がありませんし、Aさんの学籍番号はBさんの学籍番号よりも小さいのでAさんの成績はBさんの成績よりも下位であると言うようなことはありません。大小関係には意味がないのです。よって、学籍番号は名義尺度のデータと言うことになります。

順序尺度とはデータの区別が可能で順序や大小の評価も可能な尺度です。順序尺度のデータは一般的に数字を使って表現されますが、名義尺度と同様、算術計算には意味がありません²。順序や大小には意味がありますが、順位や大小の差は、その間隔が一定であるとは限らないため意味を持ちません。

順序尺度のデータの例として100m走の順位を挙げることができます。ある大会でAさんが1位でBさんは2位、Cさんは3位だったとします。このとき、Bさんは、Aさんより遅くCさんより速い言うことはわかりますす。しかし、順位の引き算を行い、(3位 - 2位) = (2位 - 1位)なのでAさんとBさんの走力の差はBさんとCさんの走力の差と同じ、とはなりません。Bさんは1位のAさんと僅差の2位だったのかも知れませんし、あるいは3i位のCさんと殆ど同タイムの2位だったのかも知れません。順位の算術計算には意味がないことがわかります。よって、順位は順序尺度と言うことになります。

間隔尺度とはデータの順序の評価が可能で差や和の評価も可能な水準です。間隔尺度のデータの取り得る数値の間隔は等間隔で、足し算の結果や引き算の結果は意味を持ちます。しかし、0(ゼロ)が「無、何もない」と言う意味ではないため掛け算の結果や割り算の結果は意味を持ちません。

間隔データの例として気温を挙げることができます。先に見たように気温の引き算は意味を持ちます。引き算の逆演算 3 である足し算の結果も意味を持ちます。しかし, 30° C(那覇市の気温)/ 15° C(札幌市の気温) = 2となるから那覇市の気温は札幌市の気温の2倍である,とはなりません。 摂氏で表していた気温を華氏 4 で表すと 15° Cは 59° F, 30° Cは 86° Fとなります。 先程の割り算を摂氏ではなく華氏で行うと 86° F / 59° F = 約1.5と摂氏での計算結果とは異なります。 摂氏を単位にすると 那覇市の気温は札幌市の気温の2倍となるのに華氏を単位にすると1.5倍となる

尺度水準

また、データは連続なのか、非連続なのかで

離散値(計数データ) - 連続量(計量データ)

離散値: 個数や枚数など数えることができる

連続量: 計る(測る、量る) ことはできるが数えることができない、測る

データの取り得る範囲に対してデータの最小間隔が十分に小さい場合は離散量であっても計量データと見做すこともあります。例えば、人数は離散値ですが、世界の人口を対象として人数の

本来、連続量のデータが計測の都合で離散的に表現されている場合は計量データとして扱う。

例えば,

最頻值: 佐藤

データは表現する性質を基準にして更に分類することができます。この基準のことを尺度と呼び、4つの水準「名義尺度」、「順序尺度」、「間隔尺度」、「比例尺度」に分類されます(尺度水準)。

^{1.} 尺度水準 https://science.sciencemag.org/content/sci/103/2684/677.full.pdf https://science.sciencemag.org/content/sci/103/2684/677.full.pdf https://science.sciencemag.org/content/sci/103/2684/677.full.pdf https://science.sciencemag.org/content/sci/103/2684/677.full.pdf https://science.sciencemag.org/content/sci/103/2684/677.full.pdf https://sciencemag.org/content/sci/103/2684/677.full.pdf https://sci/103/2684/677.full.pdf https://sci/103/2684/677.full.pdf <a href="https://sci/1048/distancemag.org/content/sci/1048/distancemag.org/co

^{2.} データに対して間隔は一定であると言う前提条件を設定することで算術計算に意味を持たせるようにすることもあります。 $\underline{\boldsymbol{e}}$

^{3.} ある演算によりAが <u>←</u>

^{4.} 華氏温度をF, 摂氏温度をCとすると\$F = \frac{9}{5} + 32\$となる €