

クロス集計と感度・特異度，データの可視化

今日は、「クロス集計」と「データの可視化」の2つのトピックを扱います。統計調査によって集められるデータには、今回説明する「尺度水準」で述べるように、質的データと量的データに大きく分けられます。大雑把に言えば、量的データとは平均（算術平均）に意味のあるデータで、質的データはそうではないデータです。次回以降の講義では、平均から出発してデータ解析の手法を説明していきますが、今回は平均ができない質的データの簡便な整理法として、クロス集計を説明します。さらに、クロス集計のしくみを使って、いま問題になっている「検査の信頼性」を、「感度」と「特異度」という概念を使って説明します。また、データをグラフによって直感的に把握できるようにする「可視化」についても説明します。

尺度水準

調査によって集めたデータは、ふつう数値で表されています。というよりも、統計学は、集めたデータに対して計算をすることで、データの集まりから情報を取り出そうとするものですから、数値で表されたデータを用いるのがふつうです。

ただ、データが数値で表されているからといって、必ずしも「数量」を表しているとは限りません。例えば、三択問題で「1番・2番・3番さあどれ？」というとき、1,2,3は選択肢の名前に過ぎず、a,b,cでもイ・ロ・ハでも同じですから、数量を表してはいません。

そこで統計学では、数値で表されたデータを、それが数量としての意味をどの程度持っているかによって、4つのレベルに分類しています。これを**尺度水準**といいます。

質的データ

名義尺度

一番レベルが低いのは、**名義尺度**です。これは、さきほどの「三択問題の1番・2番・3番」や「男性：1，女性：2」のような数値で、数値は選択肢を区別するためだけにあり、2番が1番より「大きい」という意味はありません。

順序尺度

次のレベルにあたるのが**順序尺度**です。これは、「この講義に満足しましたか？ 1) 非常に不満・2) 不満・3) 満足・4) 非常に満足」といった調査で得られる数値です。この例では、番号の順序に意味があり、2番には「満足度が1番より大きい」という意味合いがあります。しかし、「1番と2番の満足度の差」と「2番と3番の満足度の差」が同じということはありませんし、ましてや4番が2番の2倍満足しているということもありません。

名義尺度と順序尺度にあたるデータを、**質的データ**といいます。質的データは、足し算引き算をすることができません。

量的データ

一方、さらに上のレベルのデータは、足し算引き算ができるデータで、これを**量的データ**といいます。量的データは、さらに次の2つのレベルに分けられます。

間隔尺度

間隔尺度は、数値の間の順序だけでなく、数値の間隔にも意味のあるデータです。例えば、摂氏温度は間隔尺度で、「0℃と10℃の差」と「10℃と20℃の差」はどちらも10度で、同じ意味があります。しかし、20℃が10℃の2倍暖かいという意味はありません。もしそうなら、20℃は-10℃の何倍暖かいのか？ということになってしまいます。

比例尺度

間隔尺度の性質を持ち、さらに「データ間の比率」にも意味があるのが、最上位のレベルである**比例尺度**です。例えば、40歳の方は20歳の方の2倍の年数を生きていますから、年齢は比例尺度です。温度でいえば、絶対温度（それ以上冷やすことのできない「絶対零度」をゼロ度とした温度）は比例尺度で、絶対温度が2倍であれば2倍のエネルギーを表しています。

データの整理の方法として、平均（算術平均）がよく知られていますが、算術平均はデータを足し算してデータの数で割ることですから、量的データでなければ意味がありません。¹

クロス集計

クロス集計とは

次回以降の講義では、平均を計算できる量的データを対象として、記述統計学によるデータ解析の手法、さらに標本調査を用いた統計的推測の方法を説明していきます。今日は、質的データに対する解析の手法として、**クロス集計**について簡単に説明します。

例えば、「商品Aが好きか嫌いか」を調査し、「好き：50%、嫌い：50%」というデータを得たとします²。これだけでは、「好きな人と嫌いな人の割合は同じ」ということ以上の情報は得られません。そこで、この調査のさいに、回答者が男性か女性かも調べておいて、男女別にデータを整理します³。その結果を、表1のようにまとめます。

このように、ひとつのデータ群を2つの項目から見て、それらの関係を表に表すのが、クロス集計です。この表によると、男性については商品Aが好きな人がいくぶん多く、逆に女性は商品Bが好きな人がいくぶん多いことがわかります。

¹上で順序尺度の例としてあげた「授業評価」でも平均点を出していることがありますが、厳密には意味がないこととなります。ただ、このような調査では「各番号の満足度の間隔が概ね等しい」つまり近似的に間隔尺度であると仮定して、平均にも意味があるとする考え方もあります。

²ふつうは「どちらでもない」という選択肢もあると思いますが、今日は説明を簡単にするためにこのようにしておきます。

³最近では、性別を尋ねる質問には注意を要します。現在では「男性・女性・その他または答えない」という選択肢の質問が多くなりましたが、ここでは、説明を簡単にするために「男性/女性」としておきます。なお、関西大学では2017年度から、受講者名簿に性別欄がなくなりました。

表 1: クロス集計

	好き	嫌い	合計
男性	20	30	50
女性	30	20	50
合計	50	50	100

表 2: 感度・特異度

	本当に病気である	本当は病気ではない
検査で陽性	A	B
検査で陰性	C	D
合計	A + C	B + D

「感度」と「特異度」

ここで、クロス集計と同様の方法を使って、「検査の感度」についてお話ししておきます。ある病気の新しい検査法が開発されたとして、本当に病気であるかどうかわかっている人に対して、その検査法を適用して、本当にその検査法が有効かどうかを調べる実験を行います。実験結果を、「本当に病気の人」「本当は病気でない人」のそれぞれの人数と、「検査で陽性」「検査で陰性」のそれぞれの人数とで、クロス集計の形で表したものが表 2 です。

感度

検査の**感度**とは、「本当に病気の人のうち、検査で陽性になった人の割合」で、表では $\frac{A}{A+C}$ にあたります。もちろん感度は高いほうがよいのですが、それだけではその検査法が優秀だとはいえません。なぜなら、「病気の有無にかかわらず、いつも陽性と答える」検査なら、 $C=0$ ですから、感度は 100% になるからです。もちろん、こんな検査には意味はありません。

特異度

そこで、検査の能力を表すには、感度以外に**特異度**というものも用いられます。特異度とは「本当は病気でない人のうち、検査で陰性になった人の割合」で、表では $\frac{D}{B+D}$ にあたります。感度は病気の人について、特異度は病気でない人について、それぞれ検査が正解を答えている割合です。どちらも高いほうがよいのですが、両方を同時に高くすることは難しく、実際には「特異度 90% のときの感度がいくら」といった表現で、検査の能力を表します。

検査の信頼性

では、こんな問題を考えてみましょう。

ある病気を、感度 80%、特異度 99% で発見する検査があります。この病気にかかっている人が検査対象者の 1% であるとき、検査で陽性だった人のうち、本当にこの病気にかかっている人の割合はいくらでしょうか。

簡単に説明するために、検査対象者が 10000 人いるとしましょう。すると、この病気にかかっている人は、10000 人のうちの 1% だから 100 人となります。あとの 9900 人はかかっていません。

この検査は「感度 80%」ですから、病気にかかっている人 100 人のうち 80 人は、検査で陽性となります。あとの 20 人は、「病気にかかっているのに陰性」すなわち「偽陰性」となります。つまり、20 人の病気を見逃してしまうわけです。

一方、この検査は「特異度 99%」ですから、病気にかかっていない人 9900 人のうち、 $9900 \times 0.99 = 9801$ (人) は、検査で陰性となります。残りの $9900 - 9801 = 99$ (人) は、「病気にかかっていないのに陽性」

表 3: 検査の信頼性

	本当に病気	本当は病気でない	合計
検査で陽性	80	99	179
検査で陰性	20	9801	9821
合計	100	9900	10000

すなわち「偽陽性」となります。つまり、この99人については、病気でないのに治療したり隔離したり、ということになります。

以上の数値をまとめたのが、表3です。この表で、検査で陽性の人は、 $80 + 99 = 179$ (人) います。その179人の中で、本当に病気なのは80人だけです。したがって、「検査で陽性だった人のうち、本当にこの病気にかかっている人の割合」は、 $80/179 = 0.447$ 、つまり44.7%で、半分にも満たない、ということになります⁴。

「感度80%、特異度99%」という、かなり信頼性が高そうに見えますが、検査で陽性だった人のうち本当に病気の人半分もいないのでは、とても信頼できません。ではどうすればいいのでしょうか？ それは、演習問題で扱います。

データの可視化

収集したデータをまとめて図に描いてみることは、データの傾向をつかみどのような解析手法を使うかを考える、データ解析の第一歩になります。また、データを用いて人を説得するためには、グラフによる表示は説得力を飛躍的に高める道具になります。しかし、グラフを見るほうとしては、グラフから作者の意図を見抜き、だまされないようにする必要があります。今日は、皆さんがおそらくよくご存じの「グラフ」を例にとり、見た目に対する注意を考えてみます。なお、この講義では、この先にも「ヒストグラム」「散布図」という可視化の方法が出てきます。

棒グラフ

各データの大きさを棒の長さで表して比較する「棒グラフ」は、小学生の時から知っているおなじみのグラフですが、慣れてるだけに、よく注意して見ないとだまされるおそれがあります。次の設問について考えてみましょう。

1. 図1は、いずれも同じデータをグラフにしたものです。差が際立って見えるのはどれでしょうか。
2. 棒グラフの棒を、いろいろな形で表すと、親しみが持てるグラフにはなりますが、誤解を生みやすくもなります。図2の例は、棒グラフの棒を缶の形にしたものですが、このような描き方は正当でしょうか？

右のグラフは、棒の長さが値と比例しておらず、棒の長さどうしの比率を誇張することで、各棒の値の差を際立たせています。左のグラフは正しく書かれていますが、各棒で差がある部分を小さくなるので、差がないかのように見えます。真ん中のように切れ目（ブレイクといいます）を入れるのが、まあまあフェアなやりかたでしょう。

⁴この計算で用いた考え方を定理として表したものを、**ベイズの定理**といいます。

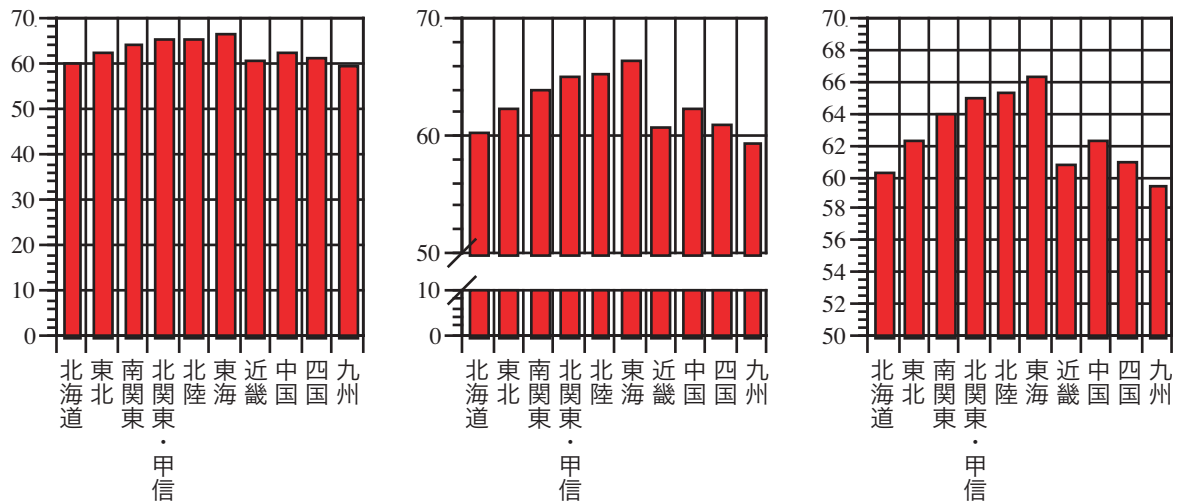


図 1: 棒グラフの例 (平成 9 年就業構造基本調査より)

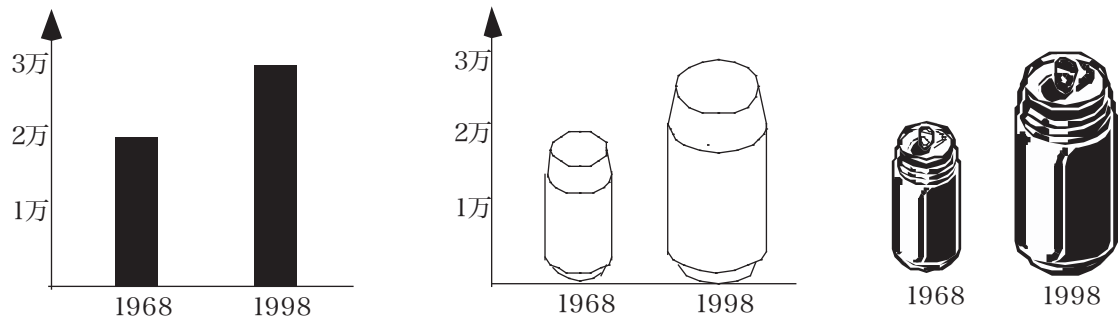


図 2: 怪しいグラフ (架空のデータ)

真ん中のグラフは、縦軸を表示して、棒の高さが値を示すことを明らかにしてはいます。しかし、棒の高さだけでなく面積も拡大されていますので、面積が値を示しているかのように見えます。高さが2倍になれば面積は4倍になりますから、より差が際立って見えるようにして、読者を欺こうとしています。

右のグラフは、体積で値を示しているかのように見せています。高さが2倍になれば体積は8倍になるので、さらに差が際立って見えます。しかも、このグラフには縦軸がないので、高さで値を示していることがわからず、体積で値を示しているかのように錯覚させています。これは完全に反則です。

ナイチンゲールのグラフ

フローレンス・ナイチンゲールは、近代看護の先駆者として知られており、その誕生日である5月12日は「看護の日」(国際ナースデー)とされています。なぜか日本では「やさしいかんごふさん」のように伝えられてきましたが、ナイチンゲールの業績は、

1. 戦場の病院での衛生管理を徹底することによって、感染症による死者を大幅に減らしたこと
2. その実績をデータとしてまとめ、グラフによる可視化を行って英国の女王や政治家に示すことで、

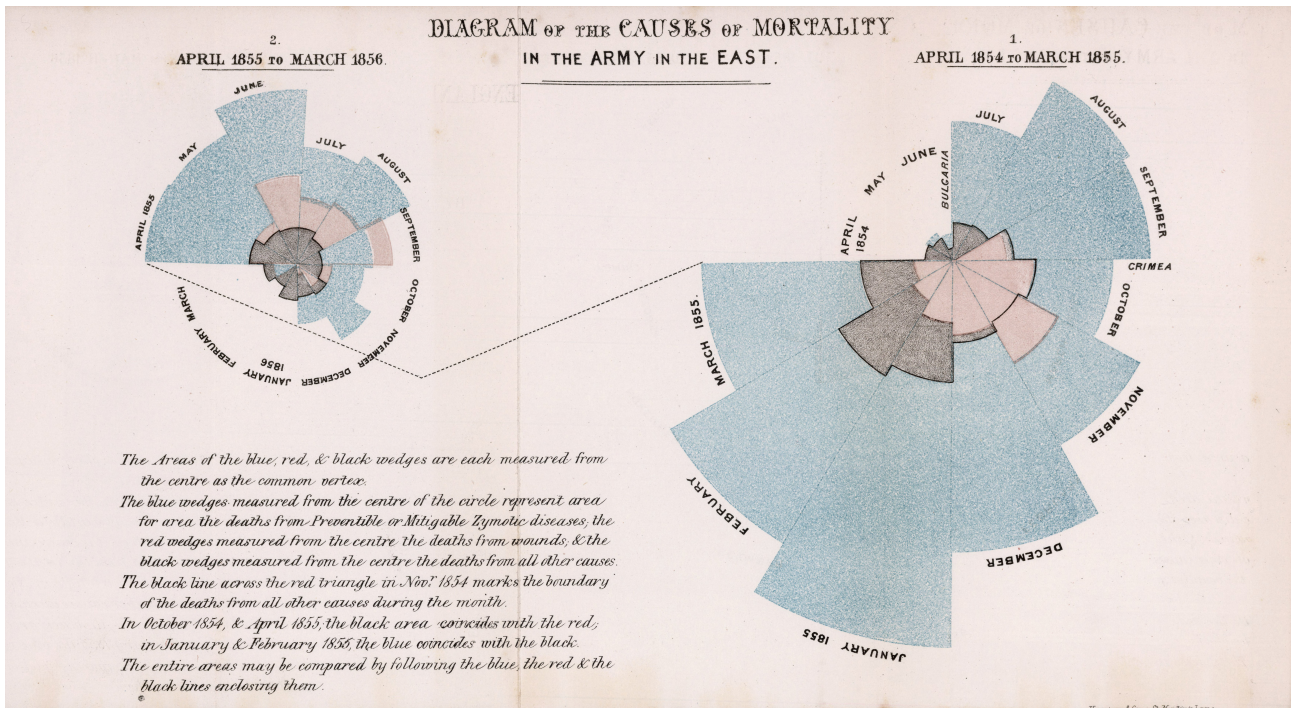


図 3: ナイチンゲールのグラフ

同国の政策を動かしたこと

にあります。

図 3 は、ナイチンゲールが作った、クリミア戦争における英国陸軍の死者数のグラフです⁵。各々の扇形の面積が、各月の死者数を表していて、色分けが原因を表しています。大きな部分を占めているのは感染症による死者で、戦傷そのものよりも感染症がおもな死因になっていること、衛生管理によってそれが大きく減少したことを示しています。

しかし、「ナイチンゲールのグラフ」は、現代の統計学の水準では、問題のある可視化です。このグラフでは、青・赤・黒の扇形の面積で死者数を表しています。面積で量を表すと、大小の印象はつきやすくなります。しかし、長さや角度に目盛りをうつことはできますが、面積に目盛りをうつことはできないので、面積で表された量を正確に把握するのはむずかしいです。とくに、正方形等でなく扇形の面積で量を把握するのはむずかしく、そのため、グラフが正確にデータを表現しているのかどうか確かめることもむずかしいということになります。

いわゆる「円グラフ」も面積で量をあらわしていますが、扇形の角度と面積が比例しているので、グラフを見て量を把握することができます。一方、ナイチンゲール方式のグラフ (coxcomb 「鶏のトサカ」といいます) では、簡単に量を把握することがむずかしく、むしろ大きな扇形が誇張される効果があります。

⁵Wikimedia Commons より、パブリックドメイン

演習問題

本文の「検査の信頼性」の例について、次の設問に答えてください。

1. 次の場合に、「検査で陽性だった人のうち、本当にこの病気にかかっている人の割合」を求めてください。
 - (a) 検査の感度が、80%ではなく90%である場合
 - (b) 検査の感度は80%で、「検査対象者のうち、この病気にかかっている人の割合」が1%ではなく10%である場合
2. 上記の計算をもとに、検査の信頼性を高めるには何が大事かを答えてください。