

Rを利用した統計解析入門

和田 哲

(北海道大学水産学部海洋生物科学科)



この資料では、左側にスライド、右側に（つまり、いま皆さんが読んでいる欄に）話し言葉での解説があります。スライドを眺めながら解説を読めば、まるで授業を受けているかのような感じで勉強することができるでしょう（そうだったらよいなあと思っています）。

問題：研究初心者が研究を始めるとき、統計解析の方法を真剣に検討すべきタイミングは以下のうち、どれでしょうか？

1. 調査、実験前：まだデータがないとき
2. 調査、実験中：データをとっているとき
3. 調査、実験後：データをとりおえたあと



勉強を始める前に、いきなりですが問題です。

皆さんが、これから、統計解析が必要な研究、例えば、動物の行動や生態について調査や実験をする研究を始めるとします。「自分のデータでどのような統計解析を使えばよいのか」など、統計解析の方法について検討すべきなのは、1から3のうち、どのタイミングだと思いますか。ひとつ選んでみてください。

正解は、「調査、実験前」です。じつは統計解析は調査方法や実験方法に強く左右されます。そして「統計解析について何も考えずに調査したり実験したりすると、統計解析できないデータをとってしまう可能性が高い」からです。だから、調査、実験前の検討が必要です。もちろん、調査・実験の最中やその後も検討する必要はありますが。

統計解析が必要な学問分野では、統計解析できないデータは「予備観察のデータ、個人的な知見のための非公式データ」としかみなされません。統計解析できないデータは「公表できないデータ」にしかたらないのです。

ダメな統計学 p.75より改変

<https://id.fnshr.info/2014/12/28/stats-done-wrong-ja-pdf/>

多くの人が、**調査や実験が終わった後で**、統計解析の方法について相談しに来る。だがそれは、助かる見込みをもって医者に相談をするような状況**ではなく**、すでに死亡した方の検死を頼むようなものだ。

統計解析が分かるヤド



はあ…えーと…。残念ながら、すでにお亡くなりになってるやど（統計解析できないデータやど）その理由は……



でも、実際には、調査や実験が終わった後で、はじめて統計解析について悩み始めて、統計解析に詳しい人に「なんとかありませんか？」と相談する人が多いです。「ダメな統計学」という本の著者が述べているように、結論は、左に記したとおりです。「ダメな統計学」は書籍にもなっていますが、書籍のもととなった和訳の文章を、ネットからダウンロードできます。

<https://id.fnshr.info/2014/12/28/stats-done-wrong-ja-pdf/>

この「ダメな統計学」には、ほかにもお役立ち情報（というか、現在の「統計学に対する典型的な誤解」など）が満載なので、ぜひ読んでみてください。

2種類の統計解析

1. 記述統計
2. 推測統計
 - ・母集団推定
 - ・仮説検定

さて、まだ「自分のデータ」をとり始めていない皆さんも、「今のうちに統計解析について勉強しよう」というモチベーションが高まったことと思います。それでは、始めましょう！

まず、一番の基本です。統計解析は記述統計と推測統計に大別できます。そして、さらに推測統計には、母集団推定と仮説検定という2つの方法があります。

論文を読んだことがある人や、学会発表を聞いたことがある人は、たとえば「オスとメスで比較したところ、t検定で有意差が認められた」などという表現を聞いたことがあるかもしれません。これらは統計解析のうちの仮説検定の結果を説明しています。なんだかカッコいい？ でも、これだけで「意味が分かった」人はいるでしょうか。

いませんね。統計に詳しい人でも、上記の表現だけでは「で、どんな性差があったの？」と尋ね直すことになります。この表現は、大事なことを伝えていないからです。仮説検定だけでは、統計解析としては不十分なのです。母集団推定もしなければならぬし、記述統計も必要です。「t検定がなにかも分からないよ！」という声も聞こえてきそうですが、仮説検定を勉強するのは最後です。順番に勉強していきましょう。

記述統計

データを整理し、その特徴を調べ、傾向をつかみ、ひと目で分かりやすいようにまとめる作業

母集団の全数調査
(例：全国の小学生、厚岸湖のアサリ、ヤドカリ実験の結果)

データ
♂、♀、♀…
2.3, 4.2, 1.5…
12, 23, 53, …
1年, 3年, 2年…

特徴把握
表現

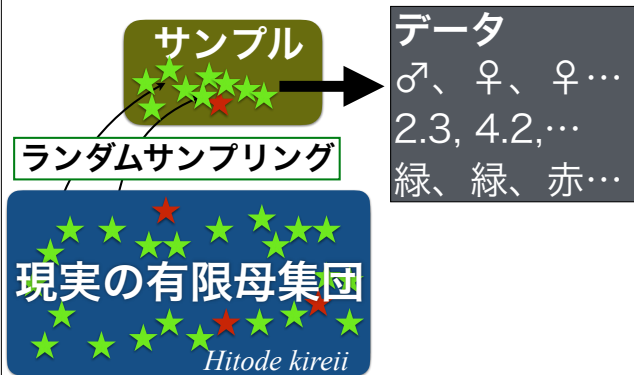
記述統計
図・表
代表値

記述統計の「記述」とは、データを記述するという意味です。つまり記述統計とはデータを表現する手法のことです。「それなら、記述統計はいつも使ってるよ。平均を計算したり、データを図示したりすることだよ」と考えた人は、きっとたくさんいると思います。正解です。

たしかに記述統計は、そのような作業のことです。でも、じつは、記述統計は、たいてい推測統計とセットになって用いられていて、そして「推測統計で不要な記述」は、すべきではありません。例えば、データの平均を計算することが不要場合もあります。そのような場合に平均を計算して提示すべきではありません。それは「記述統計においては適切でも、推測統計においては不適切」だからです。「そんなこと言われたって、推測統計ってなんのことか分からないよ」と、皆さんは思ったことでしょうか。そこで、推測統計について説明しましょう。

推測統計：母集団推定

データから母集団を推定すること



母集団推定のためには、現実の母集団からランダムに採集した標本（サンプル）が必要です。このサンプルを観察したり測定するなどして得られた文字や数値の羅列がデータです。もちろん、温度のように、現実の母集団から直接データをとることもあります。データにばらつきがないと考えられるときは、データは1つで十分であり、母集団推定は不要です。何度データをとっても結果が同じなのですから。一方、データが多少なりともばらつく（つまり、母集団から複数のデータをとったとき、データごとに文字や数値が異なる可能性がある）と考えられるときは、複数のデータをとって、そのデータのもととなる母集団はどのような集団なのか推定することになります。これが母集団推定です。

これだけでは、どのように推定すればよいのか、さっぱり分からないと思いますが、たぶん「正確な母集団推定のために、データは多ければ多いほど良い」ということは、なんとなく分かると思います。極端な話、母集団全部からデータをとれば、母集団推定しなくても良くなります。例えば、202X年の水産学部2年生の母集団における、性比（男女比）や、道内出身者と道外出身者の比率、カレーライスが好きな人の割合、視力、身長、通学時間などは、全員からデータをとれば、推定する必要はありません。

でも、科学の分野では、たいていそうはいきません。データはばらつく、けれど全てのデータを取り切ることとは不可能、だから推定しなければならない、というわけです。さて、どうやって推定するのでしょうか？

推測統計＝母集団推定

データから母集団を推定すること

- 適切な確率分布の選択
- 推定値(代表値)の解釈

データ
♂、♀、♀…
2.3, 4.2, …
緑、緑、赤…

推定

図・表
代表値

推定されるのは
無限母集団
(確率分布の母数)

ここで注意しなければならないのは、統計解析の基本では、**母集団推定は現実の母集団を推定するわけではない**、という点です。現実の母集団が有限母集団、つまり有限個のデータにすることができる母集団であっても、**実際に推定するのは、無限母集団、つまり無限個のデータで構成されていると想定された確率分布**です。確率分布は、母数というものによって、ある1つの分布に特定されます。そのため、母集団推定とは、この確率分布の母数を推定することになります。これが「統計解析が分からない」と多くの人が考える原因となる最初の(高い)ハードルです。はい、皆さんはきっと、私がなにを言っているのか分からないと思います。ですから、次回以降に詳しく説明していきますね。ただ、いまは「どうやら確率分布ってのが大事らしい」と思って、確率分布について勉強するぞと心の準備をしてください

次に、推測統計のもう1つの柱で、しかも多くの人が注目している(けれども、分かっていない人が多い)仮説検定について、おおまかに説明します。

母集団推定と仮説検定

帰無仮説とする
仮想母集団の
確率分布

検定

データ
♂、♀、♀…
2.3, 4.2, …
緑、緑、赤…

推定される母集団
確率分布

推定

まずは、母集団推定と仮説検定が、だいぶ違うものであることを強調させてください。ひょっとすると、少し統計学について学んだり、論文を読んだことのある皆さんは、t検定とか分散分析(ANOVA)とかで「有意差がある」ということを言うのが統計解析だ、と思っているかもしれません。「有意差が出れば、研究は成功したってことだ」と思っている人もいるかもしれません。このt検定とか分散分析などは仮説検定の一部です。そして、これらは母集団推定ではありません。

仮説検定は、帰無仮説と呼ばれる仮想母集団の確率分布を想定して、「データがその確率分布から得られたものだとは考えにくい」ということを主張する推測統計の手法です。統計解析の結果が「有意」であるとは、「その結果が意味の有る結果だ」という意味ではなく、「帰無仮説が正しくないように思うんだけど、、、」という意味なのです。「〇〇検定で有意な結果が得られた」などと書かれていると、まるで「判定勝ちした!」とか「実験は成功だ!」などと感じるかもしれませんが、そういう意味ではないのです。

仮説検定は、データをとってきた母集団については、なにも想定していません。そのため、t検定とか分散分析などの結果を示す人は、同時に母集団推定の結果も示さなければ、データを見たことがない人には「その有意差って、どういうことなのか」が分からないことになります。多くの論文では、その母集団推定の結果を、文章や図で伝えています。でも、それが母集団推定の結果だと分かっていない人もたくさんいます。

母集団推定と仮説検定

帰無仮説とする
仮想母集団の
確率分布

仮説検定に必要な能力
・適切な検定法の選択
・検定結果の解釈

検定

データ
♂、♀、♀…
2.3, 4.2, …
緑、緑、赤…

仮説検定は、データをとってきた母集団については何も想定していませんが、だからといって、もとの母集団を無視した検定をおこなうことは、結果を間違える原因となります。帰無仮説の確率分布が、もとの母集団の確率分布に合ったものになっている必要があります。

なるほど、母集団推定と同じように確率分布をちゃんと想定すればいいってことだね、と考えた人は、正しいのですが、それでもちょっと間違えています。じつは、例えばt検定は、たしかにt分布という確率分布を利用した検定なのですが、実際の母集団がt分布となることは（たぶん必ず）ないのです。なぜかという、t分布は仮説検定用の確率分布だからです。

このように、確率分布とひとくちに言っても、そのなかに母集団推定用の確率分布と仮説検定用の確率分布があります（統計学の入門書を読んでも、このように説明している本は見当たらないのですが、、、）。これが「統計解析が分からない」と多くの人考える原因となる2つ目の（高い）ハードルです。いまは、なにを言っているのか分からないと思いますが、先ほどと同じように「やっぱり確率分布ってやつが重要らしい」と思っておいてください。

記述統計の役割

帰無仮説の
確率分布

検定

記述 (下準備) 母集団推定や仮説検定の方針を立てるためにデータを図表にすること

データ



推定母集団の
確率分布

推定

推測統計の概要について話したので、記述統計に戻しましょう。記述統計は推測統計を支える重要なものです。私たちは、データを手に入れたら、推測統計の方針を決めるために、データの記述をおこないます。記述統計は、推測統計の下準備になるわけですが、それだけではありません。

実際におこなった母集団推定や仮説検定の手法が適切であるかどうか、そのデータのことを知らない人にも判断できるように、データを記述する必要もあります。つまり、私たちは、データの統計解析の結果を発表するためにもデータの記述をおこなう必要もあるのです。

では、記述統計について、もっと詳しく説明していきましょう。