

統計学入門

小波秀雄

Oct. 2017

はじめに

多数のデータから意味のある情報を抽出するのが統計的手法であり、その理論が統計学 (statistics) である。統計学は、確率論を基礎にして、不確実性を含む多数のデータから、一定の確実さをもった判断を下すことを目的にしている。

統計学は、社会や人間に関わるさまざまな事象の分析と多数のデータの定量的な取り扱いを可能にすることから、社会科学や医学などの人間集団を相手にした学問研究分野、心理学や教育学などの人間行動の分野、品質管理などの生産現場、保険や経営といったマネジメント分野、また政策決定のための指針作成など、さまざまな分野で広範に活用されている。

自然科学の分野でも、不確実性を含む自然現象は数多く、データの統計的な取り扱いが必要になる。また情報理論の中でも確率論とその応用は重要な一分野である。特に本書で展開される確率分布の理解は情報理論の中でも基本となるものである。

このように、統計学はまさに現代の学問と産業を支えている主要な理論のひとつであるといっても過言ではない。

その反面、確率や統計の誤った解釈や、意図的に捻じ曲げられた解釈によって、誤った指針や主張が導かれることも稀なことではない。嘘をつくための道具として、統計が不審の眼を向けられることも昔からよくあることである。誤った解釈に振り回されたり、統計の嘘にだまされたりしないためにも、統計の理論を基礎から理解することは大切である。

この講義では、確率論の入門からはじめて、古典統計学の入門的な部分を一通り取り扱う。社会学系のための統計学としては大体網羅してあるが、 F 分布など統計分布の一部を割愛した。この程度の知識があれば、実社会において現れる統計の意味をほぼ理解できるはずだが、現代的な多変量統計や予測統計については、さらにこの先の学びの課題として考えてほしい。

表計算アプリケーションの利用

統計処理では数多くのデータを使って多数回の計算を行う。

その労力を省くために、Excel や Numbers, OpenOffice/LibreOffice Calc といった表

計算アプリケーション^{*1}を使うと便利だ。セルにデータを打ち込んでから、簡単な数式を使って一斉に同型の計算をさせたり、総和を取ったりできるので、このテキストの問題を解くために活用してみることをお勧めしたい。

ただし、これらを使用する際に注意しておかなければならないのは、特に統計関数を利用したときに、出てきた数字をそのまま信用してしまって、ミスを見逃してしまうことである。たとえば分散および標準偏差を求める関数として VAR と STDEV があるが、これは第 1 章で出てくる標準偏差とは定義と値が異なることを知っておかないとまずい。

統計計算のためのパッケージの利用

本格的な統計処理のためのパッケージとして、SPSS や SAS などといったアプリケーションが知られている。特に近年では、オープンソースの統計処理のためのプログラミング言語である R^{*2} が開発されて、広く利用されるようになっているので、これから何らかの統計処理パッケージを導入する場合には、まず R を使うことをおすすめしたい。単に「R」で検索するだけでダウンロードの仕方も含めて情報が手に入るようになっている。R については、多数の参考書やマニュアルも出版されているので、その意味でも学びやすい環境になっている。巻末に R に関する情報をまとめてあるので参考にいただきたい。

正しくアプリケーションを使うために

車を運転するのに、エンジンの仕組みや道路設計に関する知識は必要ない。それでも、どこに行こうとしてハンドルやアクセルを操作しているのかを分かっていると、車はあらぬところに到着してしまう。ところが、それでも「目的地に到着しました！」と運転手が宣言する、そんなことがあったら客はどう思うだろうか？

ところが、「コンピュータで統計処理をやりました」といって、これとまったく同様の誤りを犯してしまうことはむしろありふれている。研究や実務に携わる人でさえ、実は統計学について無知なままに手続きだけを覚えて、結果を出しているケースは珍しくない。それを避けるには、統計的なデータ処理の意味をわかっておくことが必須であり、このテキストはそのために書かれている。

統計学を学ぶということは、難しい数学をマスターすることではないし、まして、基本的な定理の証明にまで遡って勉強する必要はないと言える。このテキストでも、ほとんどの数式の導出は付録に回して、数学的に納得したい人の便宜を図りながらも、本文では数学的な細部にあまり立ち入らないように留意した。

しかし、数値データを材料として処理を進める以上、その処理が何を意味しているかを

^{*1} Excel, Numbers, はそれぞれ Microsoft Office, Apple iWork, に含まれる表計算アプリケーション。

^{*2} R はフリーソフト財団の GNU プロジェクトとして開発されているので GNU R と呼ぶこともある。

理解するためには、最低限の数学的な扱いは必要である。それを押さえた上でアプリケーションの使い方をマスターすれば、安心して、かつ創造的に統計の手法を活用できる人になれるのだ。そのつもりで、本書を学んでほしい。

インターネットで利用できる問題演習システムについて

著者が開発したオンラインの問題演習システムが京都女子大学のサイトに用意されています。ゲストとしてであれば誰でも自由にアクセスして利用することができますので、本書と一緒に活用していただけると幸いです。URL は下の通りです。

<http://ruby.kyoto-wu.ac.jp/Statistics/Training/>

この本の利用について

この本の PDF ファイルは下からダウンロードできます。

<http://ruby.kyoto-wu.ac.jp/~konami/Text>

ダウンロードは自由に行っていただいてもかまいません。利用にあたっては、次の点に留意してください。

- 個人としての利用は許諾なしに行ってください。
- 学校や企業などにおける講義、セミナー等で使う際には、利用の形について著者に教えていただけると幸いです。
- 出版その他のパブリックな媒体への転載、図版の利用等については著者の許諾を得てください。
- ウェブからダウンロードできるようにするときには、古いバージョンがネット上に残ることを避けるため、上の URL へリンクすることとし、転載したファイルを別に置くことは避けてください。

著者連絡先

著者の肩書と連絡先は以下のとおりです。

京都女子大学 名誉教授

小波秀雄

E-mail: konami@kyoto-wu.ac.jp

目次

はじめに	i
第 1 章 データの整理と表現	5
1.1 データの集合から統計量を求める	5
1.2 度数分布	16
1.3 平均, メジアン, モード, どれが全体を代表するのか	19
第 2 章 初等的な確率論	23
2.1 集合と論理代数	23
2.2 集合と確率	29
2.3 認識と確率	41
第 3 章 確率分布	47
3.1 確率変数と確率関数	47
3.2 離散的な確率関数の例 — 離散型一様分布	49
3.3 離散的な確率変数の性質	49
3.4 離散的確率変数の期待値と分散	51
3.5 確率変数の関数の期待値と分散	53
第 4 章 二項分布	55
4.1 二項分布	55
4.2 多項分布	57
4.3 ポアソン分布	59
第 5 章 正規分布	63
5.1 離散的確率分布から連続的確率分布へ	63
5.2 二項分布から正規分布へ	69
5.3 正規分布表の活用	71
5.4 中心極限定理	77

第 6 章	無作為抽出と標本分布	79
6.1	無作為標本抽出	79
6.2	標本平均の分布	83
6.3	標本分散の分布	84
6.4	正規母集団	86
6.5	正規母集団と χ^2 分布	91
第 7 章	推定	97
7.1	点推定と区間推定	97
7.2	不偏推定量	101
7.3	母平均の推定	102
第 8 章	仮説と検定	111
8.1	ひょうたん島での仮説検定	111
8.2	その他の検定	124
第 9 章	相関と線形回帰	131
9.1	データの相関	131
9.2	相関係数と線形回帰	134
付録 A	重要な関係式などの導出	147
A.1	四分位数を求める	147
A.2	ベイズの定理	148
A.3	二項分布の平均と分散	148
A.4	ポアソン分布	150
A.5	標本平均の平均と分散の関係	152
A.6	標本分散の平均と母分散の関係	152
A.7	最小二乗法	153
付録 B	数表	157
B.1	正規分布のパーセント点	157
B.2	正規分布表	158
B.3	χ^2 分布表	160
B.4	Student の t-分布表	161
付録 C	ちょっとした数学的手法	163
C.1	比例配分によるデータの内挿	163
C.2	有効数字	164

C.3	数値の丸め誤差	165
C.4	多数回の計算による丸め誤差の蓄積	165
付録 D	電卓とコンピュータを活用する	167
D.1	電卓で統計計算	167
D.2	スプレッドシートで統計計算	170
D.3	本格的な統計計算には R がおすすめ	173
付録 E	解答と解説	177
索引		187

