

高額な
統計ソフトは
もういらない!?

フリーソフト

Rを使った
らくらく

医療統計解析入門

すぐに使える事例データと
実用**R**スクリプト付き

著 | **大櫛陽一** (大櫛医学情報研究所所長 / 東海大学名誉教授)

はじめに—フリーソフト R の薦めと本書の特徴—

医学をはじめ、経済学、教育学、心理学などは人間を対象とした科学・学問です。人は複雑な生き物で、同じウイルスに感染しても発病する人もいれば、発病しない人もいます。病気の治療でも、ある薬剤で大きな効果のある人もいれば、効果のない人もいます。ある人の発症する可能性や、特定の人に対する薬剤の効果を予測することは簡単ではありません。また、がんや糖尿病などの生活習慣病には、原因物質への暴露状態、栄養状態、運動習慣、遺伝的要因など多くの因子が関係しています。これらの因子のうち何が重要かを特定するには、膨大な症例の蓄積が必要で、ひとりの医師の臨床経験だけでは不十分です。

統計学を使うと、上にあげた医療上の諸問題に対しての解答（エビデンス）を得ることができます。個人ごとの疾病の発症率の予測や、特定の薬剤の治療率や副作用の発生率を予測することも可能です。統計学では、過去のデータを集めて解析しますが、その結果は目の前の健康診断の受診者や患者の近未来の生活の質を最大にするために使われるのです。

統計処理には難しい計算が必要なので、苦手だと思う学生や医療関係者が多いでしょう。しかし、最近は質の良い統計ソフトウェアが開発されていますので、計算はコンピュータに任せればいいのです。皆さんは、統計的な判断のやり方を理解して、データの収集と統計手法の選択をするだけでいいのです。この本でも、R と名付けられたソフトウェアを使って医学統計の実習ができるようになっていきます。

統計ソフトウェアには、有名な SAS や SPSS などの市販されているものもありますが、医学論文に使われているオプションを含むと 50 万円前後となり、学生個人や統計処理がときどき必要となるだけの医師にとってはあまりにも高価な商品となっています。R はフリーソフトで、インターネットの Web サイトから無料でダウンロードでき、自由に使うことができます。R の解説書は無料で提供されていますし、日本語の解説書もインターネット上に多くの研究者から提供されています。また、医学論文で使われる最新の統計処理も網羅されています。世界中の統計専門家が協力して開発し、利用者数も多いので、信頼性は市販ソフトウェアと同等です。R を使って解析された医学論文は増えており、2015 年の英語論文は 307 件ありました（PubMed で“R package”で検索）。市販ソフトウェアでは毎年のアップデートのたびに購入価格の半額程度も必要になりますから、これもかなりの経済的負担です。R は随時アップデートされており、更新もフリーなので最新の統計ソフトウェアを使えます。

本書は 2007 年に刊行した「看護・福祉・医学統計学—SPSS 入門から研究まで」（福村出版）の姉妹版です。R の出力結果は SPSS と全く同じであることを確認しています。

R で問題となるのはその操作性です。利用者としては、統計学やコンピュータの熟練者が想定されているために、簡単とはいえキーボードからスクリプトという文字列を入力して動かさなければなりません。しかも英語なので、統計やコンピュータが得意でない人では難解

かも知れません。また、利用説明書が公開されてはいますが、ものすごく多くの統計手法が掲載されており、膨大なので見るべきページを探すこと自体が簡単ではありません。さらに、統計学やコンピュータの用語で書かれていますから、やはり難解です。

本書では、医学領域で使われるほぼすべての統計手法を、簡単なものから最新の多変量解析や多重比較までを分類整理して、わかりやすく説明しています。数式は、結果を理解するために必要なものだけの最小限にとどめてあります。難解な統計処理は、PC画面の画像を掲載して、Rでの操作手順を丁寧に示しました。

本書の最も大きな特徴は、実際のデータに基づいて実習しながら理解できるようにしたこと。本書には豊富な事例を掲載しています。これらのデータは、ほとんどが実際の研究に使われたもので、統計処理の具体的な理解や、統計処理用データのつくり方に役立つと思われる。本書にある【例】では、詳細な処理手順と結果を示しており、各自のデータをつくる時に参考になり、教師が授業で説明するときの事例になります。【問】は、巻末に略解を掲示しており、読者の復習や学習達成度の確認や、教師が学生に対する課題として使うこともできます。付録の「R スクリプト一覧」や「統計処理のガイダンス」も役立つはず。

実習で必要となるR、事例データ、Rスクリプトは、中山書店のサイトからダウンロードできます。なお、ダウンロード方法については序章をご参照ください。

ダウンロードしたRスクリプトを使えば、ほとんどマウス操作だけで実習ができます。また、皆さんが収集したデータの統計処理を行いたい場合は、スクリプトに含まれるデータファイル名と変数名を置き換えれば、後はマウス操作で統計処理を行うことができます。スクリプトに慣れてくれば、本書の範囲を超えてRの大きな世界にチャレンジすることも夢ではありません。

2016年1月

東海大学名誉教授

大藪 陽一

CONTENTS

はじめに..... iii

序章 統計ソフト R のインストールと使い方

A	R の特徴と動作環境.....	1
B	R のインストール.....	1
	1. R スクリプト・R_data のダウンロード	1
	2. R のインストール	2
C	R の使い方.....	6
	1. R を走らせる	6
	2. 実習用スクリプトの読み込み	7
	3. ディレクトリ（フォルダ）の固定方法	8
	4. 実習用スクリプトを使ってみる	9
	5. スクリプトの修正と保存	11
	6. 終了手順と再開手順	11

第 1 章 統計の基礎

A	R による統計処理の基礎的知識.....	12
	1. 「EXCEL」でデータを作成する	12
	例 1.1 データの読み込みと確認.....	12
	2. R でデータセットの読み込み	13
	1) テキストデータの読み込み	13
	2) EXCEL ファイルの読み込み	13
	3. R でデータセットの確認	14
	4. R のデータセットの構造	14
	5. 変数名の明示化（データセットの 1 行目を変数名に使える）	14
	6. ケース（行）の抽出	14
	7. 一部の変数のみを抽出したデータセットを作成	15
	8. 変数の抽出	15
	9. 変数の追加	15
	10. R での注意点	16
	11. R の式と関数	16

- 12. 統計関数 16
- 13. R のデータタイプ (型) 17

B 記述統計 17

- 1. データの種類 17
 - 1) データ 17
 - ①スケール (間隔・比率尺度) に基づくもの 17 / ②順序尺度 (順位尺度) に基づくもの 17
 - ③名義尺度 (名目尺度) に基づくもの 18 / ④2 値データ 18
 - 問 1.1 次のデータはどの尺度か? 18
 - 2) 基本統計量 18
 - ①代表値 (中心傾向の測度) 18 / ②ばらつきを表す値 (ばらつきを表す測度) 19
 - コラム 数式が苦手な人のために 20
 - 例 1.2 代表値の計算 21
 - 例 1.3 ばらつきの計算 21
 - 例 1.4 基本統計量の計算 22
 - 問 1.2 23
 - 3) 分布型の話 23
 - 問 1.3 23
 - 4) 正規分布 23
 - 例 1.5 偏差値の使い方 24
 - 問 1.4 24
 - 例 1.6 中性脂肪に対する正規性の検定 (Shapiro-Wilk 検定) 25
 - 例 1.7 中性脂肪の対数変換とその正規性の検定 25
 - 5) 不偏推定量 26
 - 問 1.5 26
 - 6) 平均値の標準偏差が標準誤差となる理由 26
 - 7) グラフ表示 27
 - 例 1.8 ヒストグラムと箱ひげ図 27

C データの収集 29

- 1. 対照群の必要性 29
- 2. 無作為抽出または無作為割り当てとマッチング 30
 - ①無作為抽出 30 / ②無作為割り当て 30 / ③マッチング 30
 - 問 1.6 30
 - 問 1.7 30
- 3. 無記名アンケートと二重盲検法 31
- 4. 層別化 31
 - 例 1.9 層別化を必要とする例 31

D	統計的判断とは	32
1.	仮説検定	32
2.	両側検定と片側検定	33
3.	仮説検定の立て方と検定用統計量	34
4.	統計処理についての手順と注意	34
	①研究計画とデータの収集 34 / ②統計前処理 35 / ③統計ソフトの選択 35	
	④統計手法の選択 35 / ⑤比較の方法 35	

第2章 2群の比較

A	母集団と標本との比較	36
1.	母平均と標本平均の比較 (スケールの場合)	36
	1) 正規分布するデータで母 SD が既知の場合: Z 検定	37
	2) 正規分布するデータで母 SD が未知の場合: t 検定	37
	例 2.1 母平均と標本平均の比較 (1 標本 t 検定)	38
	問 2.1 全国平均との比較 (1 標本の t 検定)	39
2.	母比率と標本比率の比較 (1 標本カイ 2 乗検定)	39
	1) カイ 2 乗検定 (1 標本カイ 2 乗検定)	39
	例 2.2 母比率と標本比率の比較 (1 標本カイ 2 乗検定)	39
	問 2.2 全国比率との比較 (1 標本のカイ 2 乗検定)	40
B	対応のある 2 群の比較	40
1.	正規分布をしている場合: 対応のある t 検定	41
	例 2.3 2 群の比較: 正規分布をしている場合: 対応のある t 検定	42
	問 2.3 肥満対策の評価 (対応のある t 検定)	43
2.	符号付き順位和検定 (Wilcoxon の T 検定)	43
	例 2.4 符号付き順位和検定 (Wilcoxon の T 検定)	44
3.	符号検定 (S 検定: sign test)	45
	例 2.5 符号検定 (S 検定: sign test)	46
C	独立した標本の比較	48
1.	母 SD が既知で等しい場合: Z 検定 -1	49
2.	母 SD は既知であるが等しくない場合: Z 検定 -2	49
3.	等分散性の検定 (F 検定)	49
	1) 母分散は未知で F 検定で分散が等しいと判断された場合 (Student の t 検定)	50
	2) 母分散は未知で F 検定で分散が等しくないと判断された場合 (Welch の t 検定)	50
	例 2.6 F 検定の結果による t 検定	51
	問 2.4 身長 of 男女比較 (対応のない t 検定)	52

4. データが正規分布していない場合の独立した標本の比較：平均ランク検定 (Mann-Whitney の U 検定)	52
例 2.7 Mann-Whitney の U 検定	53
問 2.5 研修受講率に対する管理職の影響 (Mann-Whitney の U 検定)	55

第3章 関係を調べる

A 2変量の統計	56
1. 基本統計量	56
1) 平方和と積和	56
2) 修正済み平方和と積和	56
3) 分散と共分散	56
4) 不偏分散と不偏共分散	57
B 順序およびスケール尺度データの統計図表と 相関係数および回帰式	57
1. クロス集計表と箱ひげ図	57
例 3.1 関係を示す統計図表	57
2. 散布図と Pearson の相関係数	58
例 3.2 散布図と相関係数および回帰式	59
問 3.1 身長と体重の関係 (散布図, 相関係数)	61
3. Spearman の相関係数 (順位相関係数)	61
例 3.3 肥満度と循環器判定 (クロス集計表と Spearman の相関係数)	62
問 3.2 理解能力と表現能力の関係 (クロス集計表, 順位相関係数)	63
C 名義尺度データの統計表と検定	63
1. 対応のない2群の比率を比較する	63
例 3.4 避妊教育の性差 (2×2 表)	65
例 3.5 肥満と循環器判定の関連性 (大きな表)	66
2. 母比率との比較	68
例 3.6 性行動比率の過去との比較 (母比率との比較)	69
コラム 診断精度関係用語	70
3. 対応のある2群の比率を比較する	70
例 3.7 授業の効果 (対応のある2群の比率を比較)	71
問 3.3 授業の効果 (マクネマーのカイ2乗検定)	72
D ROC 曲線	72
例 3.8 便潜血検査の有効性 (ROC 曲線)	74
問 3.4 体調, 腫瘍マーカー, 便潜血について大腸がんの診断への有効性 (ROC 曲線)	76

第4章 生存率と危険度

A	生存率	79
	1. 生存率の計算方法：Kaplan-Meier 法	79
	2. 生存率曲線の検定	80
	例 4.1 抗がん剤の副作用抑制（生存率曲線と検定）.....	81
	問 4.1 心臓移植と生存率（生存率曲線と検定）.....	83
B	危険度	83
	1. 前向き研究：コホート調査	83
	例 4.2 喫煙と肺がん（相対リスク）.....	84
	問 4.2 喫煙と肺疾患のコホート研究（相対リスク）.....	85
	2. 後ろ向き研究：ケース・コントロール研究	85
	例 4.3 大腸がんと母親のがん既往歴（オッズ比）.....	86

第5章 多変量解析

A	多変量解析とは	88
	1. 多変量データと多変量解析	88
	2. 多変量解析の分類	88
B	重回帰分析	89
	1. 重回帰モデル	89
	2. 検 定	90
	例 5.1 生物学的年齢予測式（重回帰分析）.....	91
C	多重ロジスティック回帰分析	92
	1. 重回帰分析との違い	92
	2. 多重ロジスティック関数	93
	3. 順序データと名義データの2値化	94
	1) 順序データの2値化	94
	①ダミー変数を使う方法 94 / ②科学的根拠や経験により前半と後半に分ける方法 94	
	③中央値により前半と後半に分ける方法 94	
	2) 名義データの2値化	94
	4. 多重ロジスティック回帰	95
	例 5.2 大腸がんのリスク因子（多重ロジスティック回帰）.....	95
D	Cox 比例ハザード解析	98
	例 5.3 ライフスタイルと糖尿病発症（Cox 比例ハザード回帰）.....	99

E	判別分析	102
	1. マハラノビスの距離	102
	2. 線形判別式	103
	例 5.4 複数の検査結果からの疾病の診断（線形判別関数）	103
F	主成分分析	106
	1. 主成分とは	107
	2. 主成分分析の実際	107
	例 5.5 患者が病院を選ぶ因子（主成分分析）	107
G	因子分析	110
	1. 因子負荷量	110
	2. 共通性の推定	110
	3. 因子数の決定	111
	4. 因子軸の回転	111
	5. 因子分析の実際	111
	例 5.6 患者が病院を選ぶ因子（因子分析）	111

第6章 多群の比較

A	同時推測	114
	1. 対照群との比較とすべての対の比較	114
	2. 対応のある多群の比較と対応のない多群の比較	115
	3. なぜ個別の2群の比較の単純な繰り返しではいけないのか？	115
B	独立した多群の比較	116
	1. 一元配置分散分析	117
	例 6.1 赤血球数の年齢間比較（一元配置分散分析）	119
	2. 二元配置分散分析	121
	例 6.2 年代およびストレスレベルの違いによる女性の赤血球数の比較 （二元配置分散分析）	122
	3. Kruskal-Wallis 検定	126
	例 6.3 6つの地域間での女性の赤血球数の同時比較（Kruskal-Wallis 検定）	127
	4. Mann-Whitney の U 検定	127
C	対応のある標本の比較	128
	1. 反復測定分散分析	128
	例 6.4 3つの年度間での20歳代女性の赤血球数の対応のある比較 （反復測定分散分析）	129
	2. Friedman 検定	132

例 6.5	40 歳代女性の 3 つの年度間での BMI の対応のある比較 (反復測定 of Friedman 検定)	133
3.	Wilcoxon の T 検定	134
例 6.6	ポストホック検定 (Wilcoxon の T 検定と Bonferroni の補正)	134
4.	Cochran の Q 検定	135
例 6.7	20 歳代男性の健診総合判定変化の同時比較 (Cochran の Q 検定)	136
5.	McNemar のカイ 2 乗検定	137
例 6.8	ポストホック検定 (McNemar のカイ 2 乗検定と Bonferroni の補正)	137

第 7 章 研究計画法

A	研究の目的について	139
B	研究方法について	140
1.	ケース・シリーズ研究	141
2.	断面調査研究	141
3.	ケース・コントロール研究	141
4.	コホート研究	142
5.	自己コントロール試験研究	143
6.	無作為化試験研究	143
1)	無作為抽出	143
2)	無作為割り当て	143
7.	クロスオーバー試験研究	143
C	研究計画の不備で起こる諸問題	144
1.	脱落によるバイアス	144
2.	頻度によるバイアス	144
3.	参加意識差によるバイアス	145
4.	所属グループによるバイアス	145
5.	割り当てによるバイアス	145
D	統計的判断に必要なデータ数について	146
1.	プリテスト	146
2.	計算による必要データ数の推計	147
3.	統計パッケージを使ったシミュレーションによる必要データ数の推計	147
4.	標本数の推計支援プログラム	147
E	論文の書き方について	148

付録

1	統計処理のガイダンス	149
2	正規分布の例	154
3	算数的判断と統計学的判断	155
	(1) コインを4回トスする実験	155
	(2) 箱からボールを20回取り出す実験	156
	(3) サイコロを12回撮る実験	157
4	退院患者と入院患者の疾患統計の違い	159
略	解	161
	本書で取りあげたRスクリプト一覧	172
索	引	174

第2章

2 群の比較

医療分野では、2群を比較して差があるかどうかを知りたい場合が多々ある。たとえば、

- A という薬を服用すると、症状が改善された人のほうがされない人よりもかなり多いように思えるが、確信をもってそういえるのかどうか。
- B 療法を受けたグループと C 療法を受けたグループとでは改善の度合いが違っているようにみえるが、違うといえるのかどうか。
- ボランティア経験のある学生とない学生とでは、調査結果から障害者に対する意識が違っているように思えるが、そう言い切れるのかどうか。

といった疑問が生じてきた場合である。このような場合に、統計学的に意味のある違いなのかどうかを明らかにしようとするのが、統計学における差の検定である。

本章では、2群の比較について、**A**母集団と標本との比較、**B**対応のある2群の比較、**C**独立した標本の比較、に分けて例題をあげながら解説する。

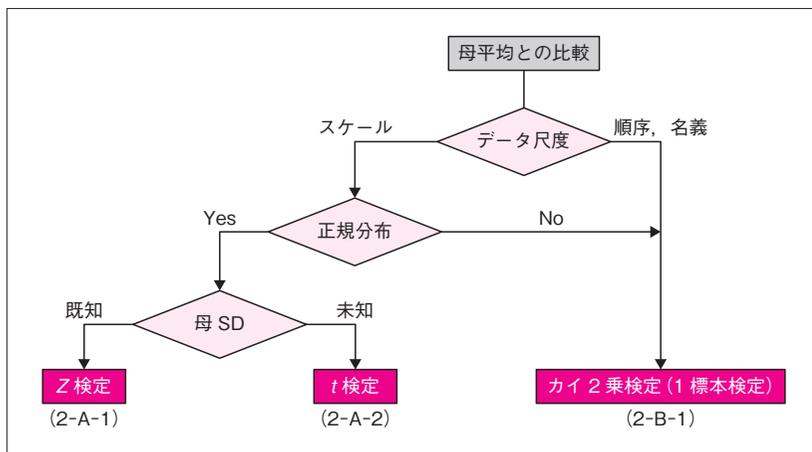
A 母集団と標本との比較

母集団と標本とを比較する場合、尺度がスケール（間隔・比率尺度）であれば母集団の平均と標本の平均とを比較することになる。また、順序尺度や名義尺度であれば母集団の比率と標本の比率とを比較する。

1. 母平均と標本平均の比較（スケールの場合）

母集団の平均と標本の平均との差を比較することがある。たとえば、14歳児の身長の全国平均（母平均）とA中学校14歳児の身長の平均（標本平均）とを比べる場合などである。比較を行う場合、データが正規分布しているのかどうか、また、正規分布をしているとしても、母標準偏差（母SD）がわかっているのかどうか、などにより手法が異なる。この手法のフローチャートを次に示す。

- 母平均
- 標本平均
- 母標準偏差：母SD



1) 正規分布するデータで母SDが既知の場合：Z検定

医療データでこの条件のそろった場合は少ない。Rでも用意されていないので電卓で計算する。

母集団の平均（母平均）、母SDがすでにわかっている場合には、次の式を用いて差があるか否かを検定する。

μ ：母平均， \bar{x} ：標本平均， σ ：母標準偏差， n ：標本のデータ数。

$$Z = \frac{(\bar{x} - \mu)}{\frac{\sigma}{\sqrt{n}}} \quad A = |Z|$$

統計仮説：標本平均（ \bar{x} ）は母平均（ μ ）と等しい。

統計仮説が成立すると $Z=0$ ($A=0$) となる。データのばらつきにより Z が 0 でないこともあるが、 Z は正規分布をするので、 Z が 0 から離れるにしたがい、その確率は大きく下がってくる。正規分布の特徴については、1章B-4「正規分布」で述べた通りである。

ここで、統計的判断は次のようになる。

Aの範囲	1.96未満	1.96以上~2.58未満	2.58以上~3.30未満	3.30以上
統計的判断	有意差なし	有意差あり	有意差あり	有意差あり
表記	NS	$p < 0.05$	$p < 0.01$	$p < 0.001$

2) 正規分布するデータで母SDが未知の場合：t検定

母集団の平均のみが既知で、母SDが未知の場合には、次の式を用いて検定を行う。ただし、 μ ：母平均， \bar{x} ：標本平均， s^* ：不偏標準偏差， ϕ ：自由度， n ：標本のデータ数である。

$$t = \frac{(\bar{x} - \mu)}{\frac{s^*}{\sqrt{n}}} \quad \text{自由度：}\phi = n - 1$$

• Z検定

• t検定

• 不偏標準偏差：
SDの不偏推定値

• 自由度：
データが自由に变化できる個数。tを計算するためにXを使うので、(n-1)個のデータは自由に变化できるが、最後の1つは決められてしまう。自由度によりt分布の形が变化し、大きくなると正規分布に近づく。

統計仮説：標本平均は母平均と等しい ($t=0$).

Rで計算をする場合には、有意確率 (p -value) をみる。有意確率の値が 0.05 未満であれば統計的に有意な差があると判断してよいことになる。

例 2.1 母平均と標本平均の比較 (1 標本 t 検定)

データは「障害施設.txt」に登録されています。

知的障害施設に入居している子ども 283 人のデータです。身長が低いように感じるのですが、文科省の全国データと比較してみましょう。身長は、年齢と性別の影響があるので、15 歳男子で全国平均身長 168.4 cm と比較します。

〈解説〉

統計仮説：平均身長は全国と一致する。

①データの読み込み

```
> dat <- read.delim("障害施設.txt")
> dat1 <- subset(dat, (年齢==15) & (性別=="m"))
> attach(dat1)
```

② Shapiro-Wilk 検定で正規性の検定を行います。統計仮説：正規分布である。

```
> shapiro.test(身長)

      Shapiro-Wilk normality test

data: 身長
W = 0.97669, p-value = 0.2948
```

• $p \geq 0.05$ なので正規性は否定されません。 t 検定が適切です。

③全国平均値と比較するために t 検定を行います。

```
> t.test(身長, mu=168.4)

      One Sample t-test

data: 身長
t = -6.7925, df = 60, p-value = 5.645e-09
alternative hypothesis: true mean is not equal to 168.4
95 percent confidence interval:
 156.1682 161.7335
sample estimates:
mean of x
 158.9508
```

- t は -6.7925 と 0 より大きく離れた値でした。自由度 (df) は 60、有意確率は 5.645×10^{-9} でした。
- $p < 0.001$ なので、統計仮説は否定され、全国と有意差があります。
- 施設の子どもの平均身長が 158.9508 cm、全国が 168.4 cm なので、約 10 cm 低いこととなります。

第3章

関係を調べる

医療分野において、データ A とデータ B との関係について調べることがしばしばある。たとえば、早起きと健康との関係、運動量と体力の関係、性別と介護に対する意識との関係など、関連性に注目が集まることはいろいろである。

データ尺度がスケールや順序で、変数 A が大きくなれば変数 B も大きくなる、その逆に変数 A が大きくなるにつれて変数 B が小さくなる、といった関係が認められる場合がある。他方、データ尺度が名義で、変数 A も B もいくつかのカテゴリーに別れ、そのクロス集計から関係が認められる場合もある。

ここでは、**A** 2変量の統計、**B** 順序およびスケール尺度データの統計図表と相関係数および回帰式、**C** 名義尺度データの統計表と検定、**D** ROC 曲線、について事例を使って学習する。

A 2変量の統計

1. 基本統計量

データ尺度がスケールの場合の基本統計量に関して、各変量の統計量に加えて、2変量による統計量の数式による表し方を示す。

1) 平方和と積和

$$\begin{aligned} S_x &= \sum x_i^2 && x \text{ の平方和} \\ S_y &= \sum y_i^2 && y \text{ の平方和} \\ S_{xy} &= \sum (x_i \times y_i) && x \text{ と } y \text{ の積和} \end{aligned}$$

- 平方和
- 積和

2) 修正済み平方和と積和

平均値を引いて計算すると修正済みとなる。

$$\begin{aligned} S_x &= \sum (x_i - \bar{x})^2 && x \text{ の修正済み平方和} \\ S_y &= \sum (y_i - \bar{y})^2 && y \text{ の修正済み平方和} \\ S_{xy} &= \sum \{(x_i - \bar{x}) \times (y_i - \bar{y})\} && x \text{ と } y \text{ の修正済み積和} \end{aligned}$$

3) 分散と共分散

$$V_x = S_x \div n \quad x \text{ の分散}$$

- 分散
- 共分散

$$V_y = S_y \div n \quad y \text{ の分散}$$

$$V_{xy} = S_{xy} \div n \quad x \text{ と } y \text{ の共分散}$$

4) 不偏分散と不偏共分散

$$V_x^* = S_x \div (n-1) \quad x \text{ の不偏分散}$$

$$V_y^* = S_y \div (n-1) \quad y \text{ の不偏分散}$$

$$V_{xy}^* = S_{xy} \div (n-1) \quad x \text{ と } y \text{ の不偏共分散}$$

- 不偏分散
- 不偏共分散

B

順序およびスケール尺度データの統計図表と相関係数および回帰式

1. クロス集計表と箱ひげ図

順序尺度または名義尺度データの2変数の関係を示すためにクロス集計表が使われる。各変量を行と列にして、表内には各行と列に対応する人数が入る。1つの変量が順序尺度または名義尺度データで、もう1つの変量がスケール尺度データのときには箱ひげ図が使われる。横軸に順序尺度または名義尺度データを取り、縦軸をスケール尺度データにする。

- クロス集計表
- 箱ひげ図

例 3.1 関係を示す統計図表

前にも使ったデータセット「意識.txt」を例にします。

〈解説〉

「介護保険が導入される前の社会福祉に対する意識」を調査する試験の結果を、年代で比較してみましょう。変数は、福祉に対する意識クイズの得点：quiz、性別：sex（1：男、2：女）、年齢3階級：age3、（1：20・30歳代、2：40・50歳代、3：60・70歳代）です。

このクイズの得点が高いほど、福祉に対する理解があるという調査です。

①データを読み込みます。

```
> dat <- read.delim("意識.txt")
> attach(dat)
```

②クロス集計表：table（行，列）

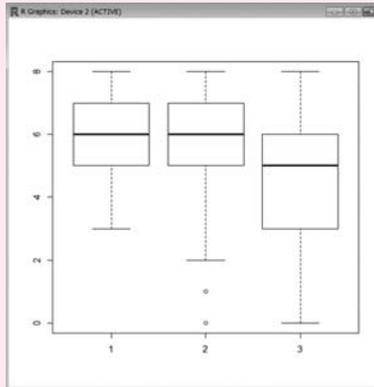
性別（sex）を行に、クイズの得点（quiz）を列にした表をつくります。

```
> table(sex, quiz)
      quiz
sex  0  1  2  3  4  5  6  7  8
  1  5  3  1  5  6 12 27 17  4
  2  8  5  3 10 14 29 32 21 11
```

③箱ひげ図：boxplot（量的変数～群分け変数）

クイズの得点（quiz）を3つの年齢階級で比較します。

```
> boxplot(quiz~age3)
```



- X軸：1：20・30歳代，2：40・50歳代，3：60・70歳代
- Y軸：福祉クイズの得点（0～8）
- 20・30歳代と40・50歳代での25%タイル値，中央値，75%タイル値は同じですが，60・70歳代ではいずれも低くなっています。この調査は介護保険導入前に行われており，高齢者では「人のお世話になる」ことへの遠慮や「他人が家に入る」ことに抵抗があったと思われます。

2. 散布図と Pearson の相関係数

2つの変数がスケール尺度データるときには，関係をみる図として散布図が使われる。原因と考えられる変数を横軸（X）に，結果と考えられる変数を縦軸（Y）にとる。この中に各データをプロットしたものが散布図で，この関係を数値で表現するのが Pearson（ピアソン）の相関係数（ r ）である。2変量の共分散を2変量の標準偏差の積で除したものであり，次の式により算出する。

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2} \times \sqrt{\sum(y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_x \cdot S_y}}$$

r が $-1 \sim 0$ のときは負（逆）相関といい， r が $0 \sim 1$ のときは正（順）相関とよぶ。 r が $+1$ に近いほど高い正の相関があるといい， -1 に近いほど高い負の相関があるという。 r が 0 であれば，2変量間に相関関係がまったくないことを意味している。

また，相関係数 r の有意差の検定は次の式で行う。

$$t = r \times \frac{\sqrt{n-2}}{\sqrt{1-r^2}} \quad \text{自由度：} n-2$$

統計仮説：相関係数（ r ） $= 0$ （ $t=0$ ）

この仮説はデータ（ X_i, Y_i ）が2次元正規分布していることを前提と

• Pearson の相関係数

中山書店の出版物に関する情報は、小社サポートページを御覧ください。
<http://www.nakayamashoten.co.jp/bookss/define/support/support.html>



大櫛 陽一（おおぐし よういち）

大櫛医学情報研究所所長・東海大学名誉教授

S22. 01. 03 生.

昭和 46（1971）年大阪大学大学院工学研究科修了，大阪府に就職。府立成人病センター，羽曳野病院，母子センター，府立病院を歴任。昭和 63 年（1988）東海大学医学部教授。平成 24 年（2012）から東海大学名誉教授，大櫛医学情報研究所所長。医療統計学，医療情報学，脳卒中，高血圧，糖尿病，メタボリックシンドローム，脂質異常症，性差医療などに関する著書多数あり。

フリーソフト R^{つか}を使った らくらく医療統計解析入門^{いりょうとうけいかいせきにゅうもん}

2016 年 3 月 7 日 初版第 1 刷発行

〔検印省略〕

著者 おおぐし よういち
大櫛陽一
発行者 平田 直
発行所 株式会社 中山書店
〒112-0006 東京都文京区小日向 4-2-6
TEL 03-3813-1100 (代表)
振替 00130-5-196565
<http://www.nakayamashoten.co.jp/>

装丁 白井弘志（公和図書デザイン室）

印刷・製本 株式会社 真興社

Published by Nakayama Shoten Co., Ltd. Printed in Japan
ISBN 978-4-521-74364-6

©Yoichi OCUSHI 2016

落丁・乱丁の場合はお取り替え致します。

- ・本書の複製権・上映権・譲渡権・公衆送信権（送信可能化権を含む）は株式会社中山書店が保有します。
- ・**JCOPY**（社）出版者著作権管理機構 委託出版物
本書の無断複写は著作権法上での例外を除き禁じられています。複写される場合は、そのつど事前に、（社）出版者著作権管理機構（電話 03-3513-6969, FAX 03-3513-6979, e-mail:info@jcopy.or.jp）の許諾を得てください。

本書をスキャン・デジタルデータ化するなどの複製を無許諾で行う行為は、著作権法上での限られた例外（「私的使用のための複製」など）を除き著作権法違反となります。なお、大学・病院・企業などにおいて、内部的に業務上使用する目的で上記の行為を行うことは、私的使用には該当せず違法です。また私的使用のためであっても、代行業者等の第三者に依頼して使用する本人以外の者が上記の行為を行うことは違法です。