

現代日本学演習 II

統計分析の基礎

田中重人 (東北大学文学部准教授)

3年生対象：2020年度 前期<金4> Google Classroom クラスコード **rugagqr****1 『講義概要』 記載内容 + α**

- ◆ 講義題目：統計分析の基礎
- ◆ 授業の目的と概要：意識調査・テスト・実験などのデータはどのように分析すればいいでしょうか。この授業では、小規模の標本調査を念頭において、統計分析の基礎的な手法を学びます。これまで統計的な分析をおこなったことのない人を対象に、初歩から講義します。同時に、コンピュータを実際に使って、データ分析の実習をおこないます。
- ◆ 到達目標: (1) 統計分析の基礎を理解する; (2) 実際にデータ分析ができるようになる
- ◇ テキスト：吉田寿夫、1998『本当にわかりやすいすぐ大切なことが書いてあるごく初歩の統計の本』北大路書房。
- ◇ 成績評価の方法：授業中の課題と宿題 (70%) と期末レポート (30%) を合計して評価する。
- ◇ その他：実習室のコンピュータ台数が限られているため、受講人数を制限することがある。
- ※ 卒業論文等で質問紙調査を予定している者は、現代日本学演習 I 「質問紙調査の基礎」(前期 金5) および現代日本学演習 V 「実践的統計分析法」(後期 金4) も受講することがのぞましい。

2 授業予定 (Google Meet 会議による)

- (1) イントロダクション (4/24)
- (2) PSPP 入門 (5/1)
- (3) 統計分析の基礎 (5/8) [序章]
- (4) 度数分布表とクロス表 (5/15-29) [1章, 4章]
- (5) 復習と進捗確認 (6/5)
- (6) 平均値の比較 (6/12-26) [2章, 5章]
- (7) 推測統計 (7/3-17) [6章, 8章]
- (8) 分析結果の表現 (7/31)
- (9) 全体のまとめとレポート内容について相談 (8/7)
- (10) 期末レポート (8/11 締切)

[] 内は、教科書の参照箇所。() 内の日付は、学期前のおおよその計画をあらわしているが、実際の授業の進行状況によって前後にずれることがある。

3 課題とフィードバック

この授業の中心は、毎回の課題とそれに対する教員からのフィードバックである。課題は <http://tsigeto.info/2020/statu/homework.html> に一括して掲載しておく。期日までに完成させて Google Classroom に提出すること。期日は、毎週**木曜日の正午**に設定してある。

- Google Classroom 「現代日本学演習 II」(LB55401 クラスコード rugagqr) に登録しておくこと。
- 課題に関する質問は随時受け付ける。Google Classroom での質問を推奨するが、電子メール (下記参照) その他の手段で質問を出してもよい。

- 毎週の講義時間 (金曜 4 講時) に、Google Meet による議論の時間を設ける。すでに提出した課題に関する質問のほか、これから先の課題や、関連する話題についての質問も受け付ける。
- 提出された課題の内容によっては、再提出を指示することがある。また、特に指示がない場合も、書き直したものを再提出してよい。
- 課題は 1 回につき 6 点。最初に提出された内容でいったん点数をつけるが、再提出された場合には加点することがありうる。
- 何を調べてもよいし、誰と相談してもよいが、それらの情報源について解答の中で説明すること。

これとは別に期末レポート (8/11 締切) がある

クロス表と平均値の比較の両方について適当な分析をして結果を解釈する。それぞれ推測統計 (区間推定または統計的検定) の結果もつけること。図・表は読みやすく整形し、論文としての体裁を整えること。授業で配布した以外のデータを使ってもよいが、その場合はデータについての解説をレポート中にふくめること。

4 受講環境について

4.1 PC 等の準備状況

この授業では、つぎの 2 種類のソフトウェアが必要である

- 統計分析
- グラフ作成

前者については PSPP、後者については Excel 等の表計算ソフトを使うことを考えている。初回授業時に受講者各自の利用環境を調査して、利用できない受講者がいる場合は変更することがありうる。

4.2 PSPP について

PSPP は、SPSS (という広く普及している有料の統計分析ソフト) に似せて作られた無料のソフトウェアです。試しにインストールしてみて、うまく動くかどうか確認して、結果をお知らせください。

- Windows の場合、<https://sourceforge.net/projects/pspp4windows/files/> から「Download Latest Version」をクリックしてください。
- 他の OS については、<https://www.gnu.org/software/pspp/get.html> を見てください

5 受講者の興味と数学的知識の調査

Google Classroom に課題を出しているのので、答えておいてください (木曜正午まで)。これは受講者の予備知識とレディネスについて知るためのもので、採点対象外です。わからない場合は「わからない」と書いておいてください。

第2講 PSPP 入門

田中重人 (東北大学文学部准教授)

[テーマ] データ配布 ; PSPP の基礎知識

1 前回課題について

予備知識の調査について解説

2 模擬データ入力実習

配布した架空の回答票 (別紙) をもとに、データを入力してみよう。

まず PSPP を起動する。通常、「データエディタ」ウインドウの「変数ビュー」タブが表示された状態になるはず。

ここで、まず「変数」を定義する。

- 変数名を必要なだけつくる。今回は a, b, ..., e とでもしておこう。変数名は自分がわかればどんなものでもよい。日本語も使える。なお、変数名以外のフィールドは入力しなくてよい
- 書き終わったら「データビュー」タブに切り替えて、いちばん上の行に変数名がなっていることを確認する。

つづいてデータを入力していく。今回は3人分のデータを用意してあって、変数は5個なので、 3×5 の行列型のデータができるはずである。

適当な名前でも保存してみる。

- PSPP データファイル (なんとか.sav) ができていることをたしかめる。
- このデータファイルは授業終了時に削除すること。(次回以降の授業ではつかわないので、おいておく必要はない。)

この方式は PSPP でデータを入力するときのいちばん簡便な方法であるが、大きなデータはあつかいにくい。実際の調査データの入力では、Excel ファイルやテキストファイルでデータを用意しておいて、PSPP に読み込むのがふつうである。

3 データ配布

この授業で使用するのは、1995年SSM調査B票の一部。調査については、配布資料のほか、『日本の階層システム』(2000年、全6巻、東京大学出版会)を参照。

- 全国から70歳以下の有権者を層化2段無作為抽出
- 訪問面接法

調査票は <<http://srdq.hus.osaka-u.ac.jp/PDF/SSM95BJ.pdf>> にもある。

ただし、配布したのはこの調査データの一部に限定したものである。

- 意識項目と基本的属性に限定 (調査票の×印はデータセットにない項目)
- 250 ケースをランダムに抽出
- 菅野剛さん (日本大学) による変数ラベルが入っている

毎回の授業で使うので、忘れないこと (調査票も)。

このデータは、この授業でのみ使用を許可されているものである。データが流出しないように注意すること。また、期末レポート提出時に、データを削除すること。

なお、自分の研究用のデータがある人は、課題などではそれを使ってもよい。ただし事前に相談すること。

4 PSPP の基礎知識

4.1 データ・セット

PSPP のデータ (「データエディタ」ウインドウで見られる) は、ケース × 変数の行列型になっている。

- 「ケース」は、個々の調査回答者にあたる
- 変数には「変数名」がついている (歴史的事情により、英数字 8 文字以内)。これだけだとわかりにくいので、変数名以外に「ラベル」をつけるのがふつう
- 無回答などの欠損値はどうなっているか?

4.2 ウインドウ構成

- データ・エディタ (上記)
- 出力ビューア (→ 分析結果やエラーメッセージなど)
- シンタックス・エディタ (プログラムを直接編集するときを使う)

4.3 分析の一般的な手続き

「データエディタ」のメニューの使いかた

- (1) 分析手法をえらぶ
- (2) 変数を指定
- (3) 必要なオプションを指定
- (4) 「OK」をクリック

結果は別ウインドウ (出力ビューア) に表示される

- 左側に目次、右側に出力内容
- エラー表示もここに出る
- PSPP のプログラム (シンタックス) も表示される

4.4 度数分布表を出してみる

- データエディタのメニュー → 「分析」 → 「記述統計量」 → 「度数分布表」
- 左側の変数リストから、分析対象とするものを選択して、右側のパレットに移動させる
- 下側の「統計」のチェックをぜんぶ外す
- 「OK」

4.5 他のアプリケーションとの連携

PSPP の出力はあまりきれいでないので、レポートを作成するときなどは、Excel や Word に表を貼り付けて整形することになる。が、出力ビューアの表をそのままコピーすることができない。

そこで、出力結果をいったん HTML で保存するなどして変換する。

- 出力ビューアのメニューから「ファイル」 → 「書き出し」を選択
- ファイル名の最後を「.html」にして保存
- 該当ファイルをブラウザで開く
- 該当部分をコピーして、他のアプリケーションに貼り付ける

アプリケーションの種類によっては、「オープンドキュメント」「CSV」等での保存も使える可能性がある。

川口秀樹 (2016) 「[PSPP]インポート、エクスポート、マージ」 <<https://note.com/xinzuzhai/n/n63b900f0bb86>>などを参照。

5 変数値の再割り当て

ウインドウ上部のメニューバーから操作する

- 「変換」 → 「他の変数への値の再割り当て」
- 変換先変数の名前をつけ、「変更」を押す。名前は英数字だけにしておくのが無難 (記号や日本語を使うと、問題がおきることがある)
- 「今までの値と新しい値」の組を順次指定する。「今までの値」は範囲で指定することも、単一の値を指定することもできる
- 値の組を指定したら「続行」を押す (元の画面に戻る)
- 「OK」ボタンを押して実行する
- 出力ビューアを右端までスクロールして、新変数ができていることを確認
- 度数分布を確認
- 問題がなければ、名前をつけてデータセットを保存 (どこに保存されるかを確認しておくこと)
- 再割り当ての手順を示したシンタックスが出力ビューアに出るので、それも保存しておくこと

6 宿題

配布したデータを使い、年齢についての度数分布表を出力する。ただし、適当な年齢幅に区切ること。結果の表を、年齢幅の設定などがわかるよう整形して、どの年齢層が多いかなどのコメントをつけて提出。また、課題の途中でどこでつまずいたかなどの経過について書いてもよい。木曜日正午までに Google Classroom に提出。

ほかの人と自由に相談してよい。

教科書のほか、つぎの資料を参考にしてよい。

- 小木曾道夫「SPSS の使い方」 <<http://www2.kokugakuin.ac.jp/~ogiso/spss/>>
- 森際孝司「SPSS の基本操作 2」 <<http://www.koka.ac.jp/morigiwa/sjs/les10201.htm>>
- 森際孝司「データ変容」 <<http://www.koka.ac.jp/morigiwa/sjs/les10401.htm>>
- 浦上昌則「SPSS おたすけマニュアル」 <http://www.ic.nanzan-u.ac.jp/~urakami/u-spss/SPSS_f.html>
- 保田時男「SPSS 操作メモ 岩井・保田（2007）準拠版」 <http://www2.itc.kansai-u.ac.jp/~tyasuda/files/2013/methoda/spss_memo_2.pdf>

ただ、これらは SPSS についての説明であるため、PSPP 操作とは一部くいちがいがある。
これら以外の資料を使ったときは、課題中に書いておくこと。

予備知識の調査：解答のポイント

(1) 1次方程式 $y = 0.5x + 1.2$ をグラフに書いたとき…

↓ ↓
傾き 切片

(2) 「必要十分条件」とは

X という条件が

あるときはかならず、そして

その時にかぎって

Y である…

(3) 「偏差値」は

平均と分散が違う複数の得点分布のなかでの相対的位置を示す

$$50 + 10 \frac{\text{生の得点} - \text{平均}}{\text{標準偏差}}$$

(4) つぎの数式の値：

$$\sum_{k=1}^{10} k = 1+2+3+4+5+6+7+8+9+10 =$$

問35 [回答票32] 次のような意見について、あなたはそう思いますか、それともそうは思いませんか。

	(ア)	(イ)	(ウ)	(エ)	
	そう思う	どちらかとう いいと思う	どちらかとう 思いはしない	そう思い はない	わからない
a 男性は外で働き、女性は家庭を守るべきである	①	2	3	4	9
b 男の子と女の子は違った育て方をすべきである	①	2	3	4	9
c 家事や育児には、男性よりも女性がむいている	1	2	③	4	9
d 専業主婦という仕事は、社会的に変意義のあることだ	1	2	3	④	9
e 女性も、自分自身の職業生活を重視した生き方をすべきだ	1	2	3	④	9
f 専業主婦は、外で働く女性よりも多くの点で恵まれている	1	②	3	4	9

問35 [回答票32] 次のような意見について、あなたはそう思いますか、それともそうは思いませんか。

	(ア) そう思う	(イ) どちらかとう どい思	(ウ) どちらかとう どい思わ	(エ) そうわ ない	(エ) わ からない
a 男性は外で働き、女性は家庭を守るべきである	1	2	3	4	9
b 男の子と女の子は違った育て方をすべきである	1	2	3	4	9
c 家事や育児には、男性よりも女性がむいている	1	2	3	4	9
d 専業主婦という仕事は、社会的に大変意義のあることだ	1	2	3	4	9
e 女性も、自分自身の職業生活を重視した生き方を するべきだ	1	2	3	4	9
f 専業主婦は、外で働く女性よりも多くの点で恵ま れている	1	2	3	4	9

問35 [回答票32] 次のような意見について、あなたはそう思いますか、それともそうは思いませんか。

	(ア) そう思う	(イ) どちらかと 思いう えはそう	(ウ) どちらかと 思いう えはそう	(エ) 思わない	(オ) わからない
a 男性は外で働き、女性は家庭を守るべきである	1	2	3	4	9
b 男の子と女の子は違った育て方をすべきである	1	2	3	4	9
c 家事や育児には、男性よりも女性がむいている	1	2	3	4	9
d 専業主婦という仕事は、社会的に大変意義のあることだ	1	2	3	4	9
e 女性も、自分自身の職業生活を重視した生き方をすべきだ	1	2	3	4	9
f 専業主婦は、外で働く女性よりも多くの点で恵まれている	1	2	3	4	9

第3講 統計分析の基礎

田中重人 (東北大学文学部准教授)

[テーマ] 度数の利用と統計分析の基礎

1 度数分布表の読みかた

- 度数
- 相対度数 (%)
- 累積度数・累積相対度数
- 欠損値のあつかい

(教科書 p. 27-31)

2 データ収集から分析まで

- (1) データの収集 (実験／観察)
- (2) 分析可能な形に加工
- (3) データ・セット作成
- (4) クリーニング
- (5) データの特徴を少数の数値に要約 = 記述統計
- (6) 誤差の評価 (この手続きの一部が推測統計)

(教科書 p. 1-6)

3 標本抽出

標本抽出の4段階モデル

- 理論母集団 = 興味の対象となる人や事物の全体
- 調査母集団 = 調査の対象とする具体的な範囲
- 計画標本 = 母集団から抽出した対象者のこと
- 有効標本 = 調査の結果あつまった有効なデータ

「無作為抽出」(random sampling) とは：

- 母集団から計画標本を選ぶ際に、母集団にふくまれるすべての個体の抽出確率が等しくなるように抽出する
- この結果として、「**確率標本**」(probability sample) がえられる

統計的な推測のための理屈は、確率標本を前提として組み立てられている。母集団の人口がわかっていて、全個体を網羅した台帳がないと、無作為抽出はできない。実際にはそういうことはないなので、いろいろ工夫して無作為抽出に近づける。

「層化2段無作為抽出」はその方法のひとつ：

- まず「地点」を抽出 (第1次抽出)
- その際、地域・都市規模等で地点抽出数を割り当てておく (層化)
- その地点の台帳から個人を抽出 (第2次抽出)

4 課題

(前回資料と同じ)

配布したデータを使い、年齢についての度数分布表を出力する。ただし、適当な年齢幅に区切ること。結果の表を、年齢幅の設定などがわかるよう整形して、どの年齢層が多いかなどのコメントをつけて提出。また、課題の途中でどこでつまずいたかなどの経過について書いてもよい。木曜日正午までに Google Classroom に提出。

ほかの人と自由に相談してよい。

教科書のほか、前回資料を参考にすること。

第4講 度数分布表とグラフの利用

田中重人 (東北大学文学部准教授)

[テーマ] グラフの種類と、それらの適切な使い分け

1 前回課題について

- 用紙上部に番号・名前を記載
- 年齢層別の人口の変動 → 1995年の人口ピラミッド <<http://www.ipss.go.jp/site-ad/TopPageData/1995.png>> と出生力の変動 <http://www.ipss.go.jp/syoushika/tohkei/Popular/P_Detail2020.asp?fname=G04-01.gif>
- 変数ラベルの利用：値の再割り当ての変数を指定する際に、変数名とは別に「ラベル」をつけることができる。あとから「データビュー」の「変数ビュー」でも
- 値ラベルの利用：「データビュー」の「変数ビュー」タブで、変数の値に「ラベル」をつける → 分析結果出力に表示される
- シンタックス (syntax) の利用

2 データセットを分割する方法

PSPPには、特定の変数の値によってデータセットを分割するコマンドがある

- メニューから「データ」→「ファイルの分割」を選ぶ
- 適当な変数を選び、「グループの比較」を選び、OK

いったんこの操作をすると、それ以降は、すべての分析が、その変数の値ごとに別々におこなわれる。元に戻すときは、「データ」→「ファイルの分割」→「すべてのケースを分析」

3 グラフの利用

分析結果は、通常、表またはグラフで示す。

表 (table): 正確な数値がわかるが、全体の傾向を読み取るには熟練が必要

グラフ (graph/chart): 全体の傾向が簡単に読み取れるが、正確さは犠牲になる

初心のうち、表とグラフの両方を作成して読んでいくのがよい

4 度数分布 (の比較) をあらわすグラフの種類

- パイチャート (円グラフ) ……半数を超えているかの判別に便利
- バーチャート (棒グラフ) ……離散量のそのままの分布を示す
- ヒストグラム ……連続量を階級に区切って示す
- 度数ポリゴン (度数多角形) ……複数の分布の比較に便利 (教科書 p. 34)
- 帯グラフ (積み上げ棒グラフ) (教科書 p. 106)

PSPP では

- 度数分布表のオプションで「チャート」を選び、「パイチャート」「バーチャート」「ヒストグラム」のなかから選択。尺度を「度数分布」または「パーセンテージ」にする。
- データエディタのメニューから「グラフ」→「バーチャート」を選び、「カテゴリ軸」に変数を指定。「バーの表示」を「ケース数」または「ケースのパーセント」に。

5 宿題 1

つぎの3種類の度数分布について、PSPP で適切なグラフを描け。

出力ビューアのメニュー「ファイル」→「書き出し」でファイル名の末尾を「.pdf」にすると PDF 形式のファイルができるので、それを木曜正午までに Google Classroom に提出。

- (1) 男女比
- (2) 本人年収の分布
- (3) 本人年収の分布の男女比較

教科書 32-37 ページを参照。

6 宿題 2

- (1) 教科書 pp. 7-16 を元に、「データの種類」の分類 (名義尺度、順序尺度、間隔尺度、比率尺度) についてまとめよ
- (2) SSM 調査の質問項目のうち、比率尺度に当たるものはどれか
- (3) などに意味があるのはどの種類のデータか
- (4) 「収入」や「学歴」を比率尺度として分析するにはどのようにすればよいか

木曜正午までに Google Classroom に提出。

第5講 クロス表分析の基礎

田中重人 (東北大学文学部准教授)

[テーマ] クロス表の書きかたと読みかた

1 前回宿題 1 について

グラフは大きく 2 種類に分かれる：

- (1) 一定の面積を分割して割合を示す：円グラフ、帯グラフ、ヒストグラム、度数ポリゴンなど
- (2) 位置または長さで量を示す：棒グラフ、折れ線グラフ、散布図など

構成比 (全部足すと 100%) を示すには (1) のグラフを使うのが原則……だが、実際にはそうでないことも多い。

- 円グラフは「半数」を基準としてみるときに使えるが、それ以外の目的には不適當
- 棒グラフ・折れ線グラフでは、縦軸の数値、目盛り、範囲に注意すること → 度数にするかパーセントにするか
- 複数の分布の比較で棒グラフを並べる場合、棒の色をはっきり違えないと、識別しにくい
- 複数の分布 (特に 3 つ以上) を比較するには、度数ポリゴンがよい (実際には折れ線グラフとして書く)
- 人数がゼロのところがあるので注意
- 欠損値のあつかい
- 連続量の度数分布は適当な階級幅に分けて「ヒストグラム」(柱状グラフ histogram) を書くのが本来である

いったん HTML ファイルに出力し、ブラウザで開いて Excel にコピーする手もある (かなりめんどくさい)

- Excel では、凡例や軸数値やタイトルのほか、各カテゴリの人数なども表示できる
- カラーで作成すると、白黒印刷では読みにくくなることが多い
- 3次元 (3D) グラフは正確な数値がつかみにくい
- ヒストグラムは Excel では描きにくい (棒グラフの距離をゼロにして、見た目をヒストグラム風にする)

2 前回宿題 2 について

- (1) 「データの種類」の分類について

- 尺度水準によって使える計算方法が違う (= 使える分析法が違う) ことを理解しておくこと
- 測定している対象そのものの性質ではなく、データにおいてどのような数値が割り当てられているかが問題である
- 上位の尺度は下位の尺度の性質を兼ね備える (たとえば間隔尺度のデータは順序尺度としても分析できる)

- (2) SSM 調査の質問項目のうち、比率尺度に当たるものはどれか → 人数、年数など

- (3) 「中央値」「四分位」などに意味があるのはどの尺度水準か → 順序尺度以上

(4) 「収入」や「学歴」を比率尺度として分析するにはどのようにすればよいか → 「変数値の再割り当て」で適当な値に変換：

- 「収入」については、各カテゴリを適当な金額に変換すればよい (たとえば各階級の真ん中の値をとる)
- 「学歴」については、その学歴を取得するのに必要な標準的年限で置き換えることが行なわれている (単に「教育年数」とよぶことが多い)。

1 → 6
2 → 8
3, 4, 5 → 11
6 → 14
7 → 17
12 → 9
13 → 12
14 → 14
15 → 16
16 → 18

3 今回の課題

「性別」と「性別による不公平」について、次の手順で「クロス表」(cross table)を作成する：

- (1) メニューから「分析」→「記述統計」→「クロス集計表」
- (2) 変数を「行」「列」にひとつずつ指定
- (3) 「セル」にパーセンテージ (行・列の両方) を追加

出力を元に、次のことを考える (参照：教科書第4章)

- この表から何がわかるか
- 「行」の%と「列」の%は何を表しているか。またこのクロス表を解釈するときにはどちらを見るのが適切か
- このクロス表をわかりやすく表示するにはどのようなグラフが適切か考え、実際に作成してみる (Excel等を使用)

第6講 連関係数とクロス表の解釈

田中重人 (東北大学文学部准教授)

[テーマ] 連関係数と%の関係を理解する

1 前回課題について

「行」と「列」については、つぎのことを理解しておくこと：

- 「行」は横方向、「列」は縦方向
- 行% と列% の使い分けは、原因→結果の関係に対応

論文等に表を載せる場合は、行%か列%どちらか一方、適切なほうだけを書く (人数は、合計のところだけ書いておけばよい)。

PSPP の結果の表を Excel 等でグラフにするには

- PSPP の出力ビューアの内容を、いったん「.html」の拡張子をつけたファイル名で書き出す
- それをブラウザで開いて、表の部分をコピー
- Excel のシートにペーストする
- 不要な部分を削除したり、ラベルをつけたりする
- 必要な部分を選択し、グラフを作成

という作業をする。

- クロス表をグラフにする場合は、帯グラフ (積み上げ棒グラフ) で合計 100%になるようにするのが標準 (折れ線グラフまたは度数ポリゴンでもよい)
- Excel の「積み上げ棒グラフ」ではカテゴリー順序が逆転するので注意 (もとどおりにしたいときは、シート上の順序をいれかえる)
- 列%によるグラフになってしまう場合は、右クリック→「データの選択」で行/列を入れ替える
- 「レイアウト」→「線」で「区分線」を指定するとよい。
- 「合計」の数値は不要

2 今回の課題

「性別」と「性別による不公平」のクロス表を作成する。ただし、つぎの数値を指定すること。

- 「セル」オプションで「度数」「期待」「残差」「標準化された残差」
- 「統計量」オプションで「カイ2乗」「Phi」

出力と教科書 (pp. 108, 116–117) をもとに、つぎのことを考える：

- 連関係数「Cramer の V」と「Pearson のカイ 2 乗」の間の数学的な関係 [式 4-19]
- 式 [4-17] のなかに、「Pearson のカイ 2 乗」「観測度数」「期待度数」「残差」「標準残差」はどのように表れているか
- 連関係数 V の最小値・最大値はそれぞれいくつか。またどのような場合に最小値・最大値をとるか。

提出は、木曜日正午まで。

なお、余力があれば、次のことも考えてみる：

- 2×2 クロス表におけるファイ係数 (ϕ : 教科書 p.110 [式 4-10]) は Cramer の V とどのような関係にあるか

3 キーワード

セル (cell): クロス表のひとつひとつの升目のこと

行 (row): 横方向のセルの並び

列 (column): 縦方向のセルの並び

独立 (無関連 = independent): すべての列について行%が等しい (またはすべての列について行%が等しい) 状態

周辺度数 (marginal frequency): クロス表の右端・下端に書く「合計」の度数

期待度数 (期待値 = expected frequency): 周辺度数を固定しておいて、独立な (架空の) クロス表をつくった場合、各セルに入る (と期待される) 度数

観測度数 (frequency): 各セルに入っている実際の度数

残差 (residual): 観測度数 - 期待度数

標準残差 (standard residual): 残差を期待度数の平方根で割ったもの

χ^2 (chi-square): 標準残差の平方和

クラメールの連関係数 V : χ^2 を全度数で割り、セル数を調整したものの平方根

行・列の数が多いクロス表では、各セルの%を比較するのが大変である。また、%の差が大きいように見えても、度数が少ない場合には、実質的には大差ないと考えるべきであるが、そのようなことを判断するのもむずかしい。そこで、まずクロス表全体について「連関係数」を見ることで、行変数と列変数の「連関の強さ」を判断し、そのうえで細かく%を比較していくのが定石になっている。

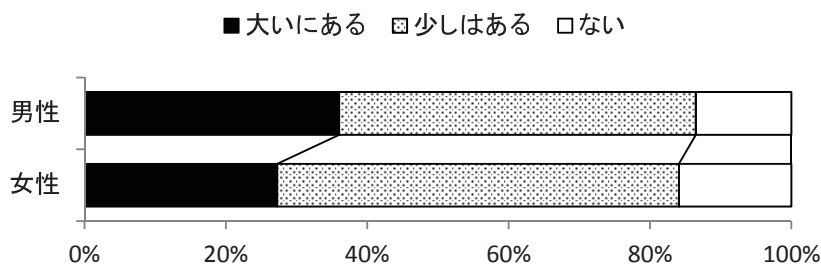


図 1 性別と性別不公平感との関連

復習と進度確認

田中重人 (東北大学文学部准教授)

[テーマ] クロス表の整形と解釈の方法

1 前回課題について

2×3 のクロス表なので、 $m=2$ である。このため、 $m-1=1$ となるので、結局 $V = \sqrt{\chi^2/N(m-1)} = \sqrt{\chi^2/N}$ である。PSPP の出力にあてはめて確認してみるとよい。

式 [4-17] のなかでは、

- 「Pearson のカイ 2 乗」 = 式全体の値
- 「観測度数」 = 「セルの度数の実際の値」
- 「期待度数」 = 「セルの度数の予測値」
- 「残差」 = 分子の括弧の中
- 「標準残差」 = 分数全体 (Σ の中身) の平方根

添字を使った表記 (たとえば n_{11} や $n_{1.}$ や $n_{.1}$ など) を覚えておくとよい。

連関係数 V は

- 独立 (無関連) のとき最小値 0
- 関連が強くなるほど大きくなる
- 完全関連のとき最大値 1 をとる

「完全な」関連とは? → 教科書 p. 115

2 連関係数の解釈

連関係数は、「モデル」と「データ」の乖離を表した値と解釈できる

- 特定のモデル (この場合、独立の状態) の下で予測される値 (この場合、期待度数) を求める
- 実際のデータの値 (この場合、観測度数) と比較する
- モデルから予測される値と実際のデータの値の違いを集計する
- 0~1 の範囲の係数になるように調整する (→ ちがうデータ間で比較しやすくなる)

多くの統計手法で、このタイプの「○○係数」が使われる。

連関係数 V の大きさの評価は主観的な問題であり、対象とする変数の性質によって基準が変わる。非常に大雑把には、つぎのような値を目安にするとよい:

- 0.1 未満 …… 関連はない (無視してよい程度)

- 0.2 程度 …… 弱い関連
- 0.3 程度 …… そこそこの関連
- それ以上 …… 強い関連

3 連関係数と ϕ 係数、%との関連

PSPP では $\sqrt{\chi^2/N}$ の値を Phi (ファイ = ϕ) として出力する。 $m=2$ のときは、 $V = |\phi|$ である。 ϕ 係数については別の計算方法がある (教科書 p. 110)。

2 × 2 クロス表においては、 ϕ 係数と%の差の間には、一定の関係がある。

- 教科書 pp. 112–114 の記述を読んで、どのような関係があるか理解する
- 実際のデータで 2 × 2 クロス表を出力して、確認してみよう

4 論文等のための表の書きかた

自分だけが見て参考にする場合はどういう表の書きかたでもよいが、人に見せる表については、きちんと整形する。

- カテゴリーの並べ順や行列の組み合わせをわかりやすく
- 変数とカテゴリーの命名
- 表のタイトルとして適切なものをつける

タイトル、表本体、注釈を読めばそれだけでわかるように書くこと。縦罫線はなるべく引かない。文字列は左揃え、数字は小数点揃えが基本。

表に書くべき要素はつぎのとおり：

- 各セルの行 (または列) %
- 行 (または列) 合計の度数と「100.0%」
- 列 (または行) 合計の%
- 全体の度数
- 連関係数 Cramer の V (または ϕ)
- 欠損数とその原因

行→列の因果を想定して表を書くのがふつうだが、ページ中の配置などの都合によっては列→行でもよい。方向は、合計の「100.0」で区別する。

- 全度数が 1000 人以下であれば、%は小数第 1 位まででよい。
- V や ϕ などの係数は小数第 3 位まで

2 列表の場合は、1 列の%だけ示してもよい

5 復習事項

5.1 PSPP の操作

- データエディタにおける「変数ビュー」の使いかた
- 「欠損値」(missing value) とは何か
- シンタックス (syntax) とは何か
- 変数値の再割り当ての方法
- グループに分割する方法
- 度数分布表における「パーセント」と「有効パーセント」のちがい
- 度数分布表における「累積パーセント」の利用法
- グラフの作成
- 「出力ビューア」の扱い

5.2 統計分析の基礎など

- 尺度水準とは何か。それはなぜ重要か。
- Excel によるグラフの書きかた

5.3 クロス表

- 「行」「列」「セル」「周辺度数」
- 「行%」と「列%」の使い分け
- クロス表をグラフにするときは、どのような種類のグラフが適切か

表 1 性別と性別による不公平感との関連

性別	性別による不公平			合計	(人)
	「大いにある」	「少しはある」	「ない」		
男性	36.0	50.5	13.5	100.0	(111)
女性	27.3	56.8	15.9	100.0	(132)
合計	31.3	53.9	14.8	100.0	(243)

Cramer's $V=0.094$. 無回答=7.**表 2** 県や市町村の部課長以上の役人に知り合いがいる比率の男女差

性別	%	(人)
男性	46.0	(113)
女性	27.6	(134)
合計	36.0	(247)

 $\phi=0.191$. 無回答=3.

第7講 平均と分散

田中重人 (東北大学文学部准教授)

[テーマ] 平均値と標準偏差の定義と計算

1 代表値と散布度

教科書 pp. 42-52 を読んで、「中央値」「四分位偏差」「平均」「標準偏差」の計算方法を理解する。

特に、表 2-1 (p. 48) で何が計算されているかを考えること。

重要なこと：

- 代表値と散布度はセットで使う
- 尺度水準による制限

2 平均値と標準偏差

平均 (mean): 総和をデータ数で割ったもの

分散 (variance): 平均値からの偏差の2乗値の平均

標準偏差 (standard deviation): 分散の平方根 (SD と書くことが多い)

教科書の表 2-1 (p. 48) を例にとると、いずれも3人分のデータセットが3つあり、

A: 2, 3, 4

B: 1, 3, 5

C: 3, 3, 3

これらのデータのそれぞれについて、つぎの計算をする

- 平均値を求める
- 「偏差」(p. 49) を求める
- 2乗する (←これはなぜか)
- 全部足す (「平方和」(sum of squares: SS) という)
- 合計人数 (=3) で割る (これが「分散」)
- 平方根をとる (これが「標準偏差」)

3 宿題

教科書 p. 52 の練習問題 2-3 について、平均値と標準偏差を計算せよ。計算の途中経過がわかるように解答すること。

第8講 平均値の比較

田中重人 (東北大学文学部准教授)

[テーマ] ふたつのグループ間での平均値の比較

1 度数分布表のオプション

度数分布表の「統計」オプションで「平均」と「標準偏差」をチェック。

- 「記述統計量」→「記述統計量」でも出力できる。(そのほかにもさまざまな分析の「記述統計」オプションで出力できる)
- PSPP などの統計ソフトの「標準偏差」は、統計学の普通の定義とはちがう式で計算している (教科書 p. 48 注6)。データが大きくなれば (およそ 200 以上なら) このことによるちがいはほとんどなくなるが、小さいデータ (たとえば 10 人程度) では大きなちがいになるので注意。

練習問題：前回宿題について、PSPP にデータを打ち込み、平均値と標準偏差を出力してみよう。

2 順序尺度の変数の「平均値」

平均値は、本来は、間隔尺度以上の水準の変数にしか使えない。しかし、実際には、一定条件を満たせば、順序尺度についても平均値をとっていいとする基準が使われている。

- 潜在的には間隔尺度のはず
- 測定のポイントが一定間隔

具体的には、4点以上の尺度であって、正規分布に近似している場合 (教科書 p. 53-59)。これは、「偶然の積み重ねで形成されるものは正規分布にしたがう」という仮定による。

「正規分布に近似」しているかどうかは、通常、つぎの3点で判断する。

- 単峰性
- 左右対称性 (歪度)
- 中央への集中度 (尖度)

PSPP 度数分布表 (または棒グラフ) で検討するとよい。

「度数分布表」の「統計量」オプションで「歪度」「尖度」を指定すると、正規分布との乖離度を統計的に検討できる。これらの値は、正規分布のとき0をとり、絶対値が大きくなるほど、正規分布から外れる。およそ ± 2 の範囲を超えていれば、正規分布からのずれが無視できない。

練習問題： 配布しているデータ中のさまざまな意識変数について、度数分布表 (棒グラフ、歪度、尖度) を出力し、正規分布に近似しているかどうかを判断してみよう。

これらの条件を満たさない場合は非線形変換 (教科書 p.142-144) をおこなったり、順位に変換したりすることがある。あるいは、平均値を使わずに中央値を使って分析することもある。

なお、2 値の変数は、この条件にかかわらず間隔尺度とみなしてよいが、一定以上のデータ数があり、あまり偏っていないことが必要。

3 平均値の欠点

平均値は「はずれ値」(outlier) の影響を受けやすい。あまりにかけはなれたケースがあるときは

- 上下数%を取りのぞく (調整平均: 教科書 p. 46)
- 順位に変換したり中央値を使って分析

などの方法を使うことがある。

また、極端なはずれ値がなくとも、左右非対称の分布の変数 (所得、人口、めったに起こらない現象の経験回数など) では、平均値より中央値の方が適切な代表値であることが多い。

4 ふたつのグループ間での平均値の比較

データをグループに分けて、それぞれ平均値 (= 層別平均) を求め、それらの間の差をもとめる。この差の大きさを、標準偏差を基準にして評価する。具体的には、effect size (ES) または 相関比 (η : イータ) という統計量を使う。

4.1 PSPP でのグループ別平均値の出力

メニューの「分析」から「平均の比較」→「グループの平均」を開く。

- 「従属変数」に平均値を計算する変数を指定
- 「独立変数」にグループの変数を指定

4.2 エフェクト・サイズ

Effect size (ES): 一般には「Cohen の d 」と呼ばれる。

$$ES = \frac{\text{グループ別平均の差}}{\text{併合 SD}} \quad (1)$$

「併合 SD」の計算については教科書 p. 137 [式 5-8] の分母を参照。大雑把には、グループ別の SD の中間の値と考えてよい。

ES は、計算が簡単であり、直感的に把握しやすい。しかし、各グループの人数を考慮せず平均値だけ比較するため、グループの人数が大きくちがう場合でも、同じ人数に 2 等分されている場合でも、その間のちがいは ES の値に反映しない。また、2 グループ間の比較だけを行うものであるため、3 つ以上のグループを比較するのにはつかえない。

ES は PSPP では計算できない。

4.3 相関比 (correlation ratio)

- 各グループの個体が全員そのグループの平均値を持つ状況を仮定してSDを求める。
- この仮想SDを実際のSDで割った数値が「相関比」である。数式では η (eta) であらわす

イータ (η) は 0~1 の範囲の値をとり、独立変数の影響力をあらわす

PSPP では、クロス表の「統計」オプションで「イータ」を選択する。出力には2種類の値が出てくるが、**平均値を計算したい変数が「従属」になっているほう**を見る。

5 課題

- (1) 適当な変数について、度数分布表・平均・標準偏差を出力（全体と男女別）
- (2) 性別でわけて度数分布をグラフに表す（度数ポリゴンまたは折れ線グラフ）
- (3) これらの分析結果から何が言えるか、解釈を書く。

第9講 分散分析

田中重人 (東北大学文学部准教授)

[テーマ] 分散分析 (ANOVA) の考えかたと計算方法を理解する

1 前回課題について

- 選択肢が4つ以上の項目を選ぶこと
- 正規分布に近似しているか確認すること (単峰性、偏り、集中度)
- グラフの縦軸の単位

2 分散分析の考えかた

グループ別の平均値を当てはめて仮想の分散を求める分析法を「分散分析」(ANOVA: ANalysis Of VAriance) という。

- 従属変数 (dependent variable) と独立変数 (independent variable)

相関比 (イータ) の性質:

- 最小値:
- 最大値:

大きさの評価基準は、Cramer の連関係数 V と同様。

なぜ相関比を求めると、平均値を比較していることになるのか?

3 課題

次のデータ (10人) について、分散分析を行なう

男性: 1, 2, 3, 3, 4

女性: 2, 3, 4, 4, 5

まず手計算 (または Excel) で考えてみて、そのあと、PSPP にデータを入力して検算する。

- (1) 全体の平均値と SD を求める
- (2) 男女別の平均値を求める
- (3) 男性の平均値 \times 5人と女性の平均値 \times 5人からなる仮想データを考えて SD を求める
- (4) (3) の SD を (1) の SD で割ったものが相関比 η

この相関比がなぜ「平均値の比較」の指標になるかを考えること。

4 相関比とエフェクトサイズの関係

相関比 η とエフェクトサイズ ES の間にはつぎの関係がある (n_1, n_2 は各グループの度数、 $N = n_1 + n_2$ は全体の度数)。

$$ES^2 = \frac{\eta^2}{1 - \eta^2} \times \frac{N^2}{n_1 n_2} \quad (1)$$

特に、2 グループの度数が等しい ($n_1 = n_2$) なら、この式は次のようになる。

$$ES^2 = \frac{4\eta^2}{1 - \eta^2} \quad (2)$$

(グループの度数が違えば、ES はこれより大きくなる)

さらに、 η があまり大きくない ($\eta < 0.4$ 程度) 場合であれば、次のような単純な式で近似できる：

$$ES = 2 \eta$$

5 モデルとデータの乖離

相関比 η は、モデルとデータの乖離を表した値と解釈できる

- 「モデル」は何か？
- データとの乖離はどうやって計算しているか？
- 係数の取りうる値の範囲は？

6 表の書きかた

- 各層と全体の平均値と標準偏差 (測定水準の2桁下まで)
- 各層と全体の人数
- 相関比またはエフェクトサイズ (小数第3位まで)
- 欠損数とその原因

7 グラフの書きかた

平均値をプロットし、上下にSDを表示する。誤差範囲 (error bar; 別名「ヒゲ」) にはSD以外を書く場合もあるので、必ず「±標準偏差」であることを明記する。

Excel では

- PSPP 出力を (HTML 形式でファイルに書き出して) シートにはりつける
- 折れ線グラフを描く
- メニューの「レイアウト」から「誤差範囲」→「その他の誤差範囲オプション」をえらぶ
- 「ユーザ設定」→「値の指定」
- 「正の誤差の値」「負の誤差の値」にSDが入っているセル範囲を指定 (おなじものでよい)

より詳細に分布の違いを検討したいときは、グループ別に度数ポリゴン (または折れ線グラフ) を描いてもよい。

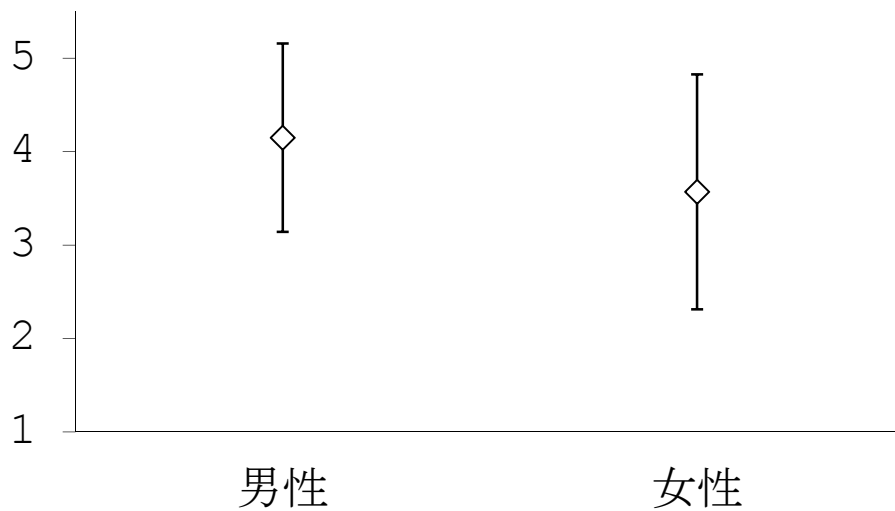
8 文献

- 奥村 晴彦 (2016) 「効果量, Cohen's d, 検出力, 検出限界」 <[https://oku.edu.mie-u.ac.jp/~okumura/stat/ effectsize.html](https://oku.edu.mie-u.ac.jp/~okumura/stat/effectsize.html)>

表 1 保守的意識の男女差

	平均	標準偏差	(人)
男性	4.15	1.01	(109)
女性	3.57	1.26	(130)
合計	3.83	1.18	(239)

「以前からなされていたやり方を守ることが、最上の結果を生む」
に対する回答: 「1. そう思う」～「5. そう思わない」
相関比 $\eta=0.244$. 無回答=11.



「以前からなされていたやり方を守ることが、最上の結果を生む」
に対する回答: 「1. そう思う」～「5. そう思わない」
相関比 $\eta=0.244$. N=239. 無回答=11.

図 1 保守的意識の男女差 (平均±標準偏差)

第10講 推測統計の基礎と区間推定

田中重人 (東北大学文学部准教授)

[テーマ] 推測統計の基礎

1 前回課題について

男性: 1, 2, 3, 3, 4 → 平均 2.6

女性: 2, 3, 4, 4, 5 → 平均 3.6

全体の平均 (SD): 3.1 (1.14)

グループ別平均値を当てはめた「仮想」データの平方和は、つぎのようになる。下線部に注意。

$$\text{仮想の平方和} = 5(\underline{2.6 - 3.1})^2 + 5(\underline{3.6 - 3.1})^2 = 2.5 \quad (1)$$

これを $N (=10)$ で割って平方根をとると標準偏差が得られる。

$$\text{仮想SD} = \sqrt{\frac{2.5}{10}} = 0.5 \quad (2)$$

$$\eta = \frac{\text{仮想SD}}{\text{実際のSD}} = \frac{0.5}{1.14} = 0.44 \quad (3)$$

ただし、SPSS では平方和を $N - 1 (=9)$ で割って「標準偏差」を求めているので、注意。度数がある程度大きくなれば (およそ $N > 200$ の場合)、このことによる違いは気にしなくてよい。

分散分析の実際の計算では、平方和どうして割り算して η を求める (N で割らずに済み、平方根を求めるのも一度で済むため)。結果を示す「分散分析表」は、伝統的にはそのようにして作成する。

2 基礎知識

- 記述統計と推測統計 (教科書 pp. 3-5)
- 母集団と標本 (第3講資料)
- 無作為抽出と確率標本

3 統計的推測のふたつの方法

- 袋のなかに色つきの玉がたくさん入っている。ここから8個取り出したところ、すべて赤であった。
- 全世界から8人を無作為抽出して麺類の好みをきいたところ、全員が「うどんが好き」と答えた。

このような情報 (= 標本統計量) から、母集団における統計量 (= 母比率) を推測する

区間推定: ある統計量の母集団における値について、確率的な推測によって範囲を求める →母比率はたぶん ○
○ から ×× の範囲にある

統計的検定: ある統計量の母集団における値について何らかの「帰無仮説」(null hypothesis) を設け、それが棄却できるかを判定する →母比率が0.5だと考えてよいか?

統計的検定のほうが計算が簡単であるため、よくつかわれている。区間推定を論文等で目にする機会はあまりないが、きちんと理解するにはまず区間推定の考え方をおさえるのがよい。

4 母比率の区間推定

4.1 区間推定の原理

- (1) 「信頼率」を決めておく (たとえば95%)
- (2) (1 - 信頼率) の確率を両極の事象に設定する (高いほう、低いほうからそれぞれ2.5%ずつを除く)
- (3) 母集団における値がいくつであれば、この両極端を除いた区間に測定値が入るかを計算する。統計量の性質に応じて、これを計算するための式と数表があるので、それを利用する。

このようにして求めた、母集団においてありうる値の集合が「信頼区間」である。通常、最初に決めた「信頼率」を明示して、「95%信頼区間」などのようにいう。

4.2 母比率の区間推定

標本の規模 n がじゅうぶん大きく ($n > 30$)、比率 m があまり偏っていない ($0.1 < m < 0.9$) とき、母比率の95%信頼区間は次の式で求められる:

$$m \pm 1.96 \sqrt{\frac{m(1-m)}{n}} \quad (4)$$

4.3 課題

全世界から400人を無作為抽出してある意見に対する賛否を聞いたところ、「賛成」と答えた人が240人であった(欠損値はないものとする)。このとき、母集団(全世界の人々)における賛成の比率の95%信頼区間を求めよ。

5 宿題

教科書 pp. 156-162 を読み、統計的検定の手続きをまとめよ

6 期末レポート

期限: 8/11 (火) 17:00

提出先: Google Classroom 「期末レポート」にファイルを提出

内容: クロス表と平均値の比較の両方について適当な分析をして結果を解釈する。それぞれ推測統計(区間推定または統計的検定)の結果もつけること。図・表は読みやすく整形し、論文としての体裁を整えること。授業で配布した以外のデータを使ってもよいが、その場合はデータについての解説をレポート中にふくめること。

備考: レポート提出後に、データのコピーをすべて消去すること。

第11講 統計的検定

田中重人 (東北大学文学部准教授)

[テーマ] 平均値の区間推定と統計的検定の方法

1 前回課題について

1.1 課題 1: 母比率の区間推定

- 母比率の区間推定においては、95%信頼区間は、 $n=100$ で±10%、 $n=400$ で±5%程度
- 母集団の規模は関係ない (無限母集団の仮定)

1.2 課題 2: 統計的検定の手続き

- 背理法的思考 (「帰無仮説」とは)
- 「臨界値」はどうやって計算するか (→教科書末尾の数表)
- 「有意でない」ことの意味
- 区間推定との関係 (「有意水準」と「信頼率」「危険率」)

2 母平均の区間推定

間隔尺度以上の変数の場合には、「母集団においては正規分布している」という仮定を置けば、平均値の区間推定が可能。標本における平均 m と標準偏差 SD から、母集団における平均 M を推測する。

95%信頼区間は次のようになる：

$$m \pm \text{臨界値} \frac{SD}{\sqrt{n}} \quad (1)$$

臨界値は、 t 分布を使って求める (数表で調べる)。「自由度」($df = n - 1$) と危険率 ($= 1 - \text{信頼率}$) によって変化する。標本規模 200 以上で信頼率 95% なら、臨界値は 1.96 と考えてよい。

3 平均値の差の区間推定

ふたつのグループの間の平均値を比較するときは、平均値のグループ間の差についての信頼区間を直接求める方法をとる。標本における 2 グループ間の平均値の差を d とすると、95% 信頼区間は

$$d \pm \text{臨界値} \times \text{併合 SD} \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad (2)$$

ただし n_1, n_2 はそれぞれのグループの人数。「臨界値」は自由度 ($n_1 + n_2 - 2$) の t 分布にしたがって求める。

4 PSPP コマンド

4.1 母平均の区間推定

「分析」 → 「記述統計量」 → 「探索的」

- 「従属変数」を指定
- 「統計」ボタンから「記述統計量」をチェック

「因子」を指定すると、グループ別に分析できる。

4.2 平均値の差の区間推定

「平均の比較」 → 「独立したサンプルの t 検定」

- 「グループ化変数」は、数値を指定しないといけない
- 連続量を一定の値 (分割点) でわけることもできる
- 出力は「独立サンプルの検定」の1行目「等分散を仮定する」を見る (この場合、「母集団で正規分布」「2層間でSDが等しい」ということが前提になる)

5 統計的検定 (statistical test)

特定の値 x (0 にすることが多い) を設定して、その値が信頼区間に含まれているかどうかを判定する。

5.1 統計的検定用語 (教科書 pp. 156–158, 165–166)

帰無仮説 (null hypothesis): 母集団における統計量が「特定の値」に等しい、という仮説

有意 (significant): 「特定の値」が信頼区間に入っていないことをあらわす

5.2 平均値の差の検定の場合

「5%水準で有意」とは……

- 95%信頼区間が x をふくまない
- すくなくとも 95 % の確率で、母集団において平均値の差があるといえる

「5%水準で非有意」とは……

- 95%信頼区間が x をふくむ
- 母集団においては平均値の差はないかもしれない

5.3 有意確率とは

信頼区間の幅は、危険率 (= $1 - \text{信頼率}$) を下げると広がる。危険率を下げても信頼区間をひろげていくと、どこかで x をふくむようになる。このときの危険率のことを「有意確率」または「 p 値」という。

分析の際は、前もって危険率を設定しておき (通常は 5 %)、有意確率がその値を下回っているかどうか 判別する。

- 有意確率が 0.007 → 5%水準で有意
- 有意確率が 0.023 → 5%水準で有意
- 有意確率が 0.088 → 5%水準で非有意

6 区間推定と統計的検定

区間推定と統計的検定の間に本質的なちがいはない。ただし、区間推定は、統計量によっては、すごくむずかしい場合がある。統計的検定のほうが計算が簡単なので、統計的検定を使うことが多い(分野によってちがう)。

7 課題

適当な変数の平均の男女間の差について統計的検定を行い、結果にコメントをつけて提出

参考資料

森際孝司「標準正規分布表」 <https://www.koka.ac.jp/morigiwa/sjs/standard_normal_distribution.htm>

第12講 さまざまな検定手法

田中重人 (東北大学文学部准教授)

[テーマ] 相関比と連関係数の検定 (F 検定、カイ 2 乗検定)

1 信頼区間と有意確率について補足

「標準誤差の差分」を 1.96 倍すると、95%信頼区間の幅の半分になる (ケース数が 200 以下の場合や、95%以外の信頼率の場合は、1.96 の代わりに t 分布表から求めた臨界値を使う)。

2 分散分析と F 検定

帰無仮説: 母集団においてはすべてのグループの平均値が等しい ($\eta = 0$)

PSPP では、

- 「平均の比較」→「一元配置分散分析」
- 「従属変数リスト」と「因子」を指定
- 「統計」の「記述統計量」をチェック

出力の「分散分析」表の右端の「有意水準」を見る。

2 グループの比較なら、平均値の差の t 検定と同じ結果。

必要とする前提も t 検定と同様 (母集団では正規分布しており、SD が全グループで等しい)。

3 クロス表の「独立性の検定」

帰無仮説: 母集団においては $V=0$

PSPP では、「クロス集計表」の「統計」で「カイ 2 乗」を指定。出力の「ピアソンのカイ 2 乗」の列の右端の「近似的有意水準 (両側)」の値を見る (各セルの期待度数が 5 以上であることを前提とする)

2 × 2 クロス表では、 χ^2 の値が大きめに出る (= 有意になりやすい) ため、種々の調整を要求されることがある。

4 課題

クロス表の「独立性の検定」と分散分析を、それぞれ適当な変数について行い、有意確率が 0.05 未満になるものを探す。その時の連関係数 V と相関比 η の値を確認すること。

第13講 分析結果の表現

田中重人 (東北大学文学部准教授)

[テーマ] 検定結果の表記

1 前回課題について

- 平均値を計算していい変数かどうか、尺度水準を確認すること
- クロス表の独立性の検定の場合、各セルの期待度数が小さいと正確な結果にならない。期待度数が5未満のセルがあるとまずいので、チェックしておくとうい。

2 検定結果の書きかた

- 検定の結果は表の下端の注釈に書く
- 検定の対象になる統計量を必ず書く
- 有意である場合は、 $p < 0.05$ のように書くか、統計量右肩にアスタリスク (*) をつける
- 有意でなければ、 $p > 0.05$ のように書くか、統計量右肩に ns と書く (= not significant)

具体例は 配布資料 PDF 版 を参照。

表そのものの書きかたについては、6月5日資料 と 6月26日資料 とを参照。

3 今後の予定

8/7: 講義はなし。質問がある場合、授業時間に来てもらえれば対応します

8/11: 期末レポート締切

表のサンプル

表 1 性別と性別による不公平感との関連

性別	性別による不公平			合計	(人)
	「大いにある」	「少しはある」	「ない」		
男性	36.0	50.5	13.5	100.0	(111)
女性	27.3	56.8	15.9	100.0	(132)
合計	31.3	53.9	14.8	100.0	(243)

Cramer's $V = 0.094$. $p < 0.05$ 無回答=7.

表 2 県や市町村の部課長以上の役人に知り合いがいる比率の男女差

性別	%	(人)
男性	46.0	(113)
女性	27.6	(134)
合計	36.0	(247)

Cramer's $V = 0.191^*$. 無回答=3.

*: 5%水準で有意.

表3 生活全般満足度の男女差 (1)

性別	平均	標準偏差	(人)
男性	2.62	1.02	(114)
女性	2.24	0.91	(136)
合計	2.41	0.98	(250)

$\eta = 0.198$. $p < 0.05$.

表4 生活全般満足度の男女差 (2)

性別	平均	標準偏差	(人)
男性	2.62	1.02	(114)
女性	2.24	0.91	(136)
合計	2.41	0.98	(250)

$\eta = 0.198^*$. *: 5%水準で有意.

表5 性別役割意識の男女差 (1)

	平均	標準偏差	(人)
男性	1.77	0.67	(111)
女性	1.89	0.65	(132)
合計	1.84	0.66	(243)

$\eta = 0.086$. $p > 0.05$. 無回答 = 7.

表6 性別役割意識の男女差 (2)

	平均	標準偏差	(人)
男性	1.77	0.67	(111)
女性	1.89	0.65	(132)
合計	1.84	0.66	(243)

$\eta = 0.086^{ns}$. ns: 5%水準で非有意.

無回答 = 7.