

現代日本論演習

3年生対象: 2009年度前期 (5セメスター: 授業コード=L52407)
<火4> コンピュータ実習室 (文学部本館 7F 711-2)

『講義概要』 p. 163 記載内容

授業内容: 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、毎回データ分析の実習をおこないます。

成績評価の方法: 各回の授業中の課題 (50%)、中間試験 (20%)、期末レポート (30%) を合計して評価する。

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。

卒業論文等で質問紙調査を予定している者は、6セメスタ開講の「現代日本論演習: 質問紙法の基礎を学ぶ」(水2) および「現代日本論演習: 応用統計分析」(木2: 大学院と合同) も受講することがのぞましい。

授業の概要

目次

1. イントロダクション (4/14)
2. SPSS 入門 (4/21)
3. 統計分析の基礎 (4/28)
4. 記述統計(1): 度数分布とクロス表 (5/19~6/2)
5. 中間試験 (6/9)
6. 記述統計(2): 平均値の比較 (6/16~6/30)
7. 推測統計 (7/7~7/21)
8. 期末レポート (8月中旬提出)

() 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。

1. イントロダクション

- この授業の概要・スケジュール・評価方法
- 部屋とコンピュータの使いかた
- SPSS の起動
- データ行列 (データセット)

2. データ配布・SPSS 入門

- データの配布
- SPSS の概要
- SPSS コマンド・シンタックス
- メニューによるシンタックス作成
- 変数値の再割り当て
- 他のソフトウェアについて (電卓, Excel, Word?)
- 印刷

3. 統計分析の基礎

- 実験と観察
- データの記述
- データの種類

4. 記述統計 (1): 度数分布とクロス表

4.1. 度数分布表

- frequencies コマンド
- 相対度数 (パーセンテージ)
- 棒グラフ・ヒストグラム・度数ポリゴン
- Excel で整形, グラフ作成

4.2. クロス表

- 度数分布表のグループ化
- クロス表表記
- 行と列の%
- 周辺度数 (marginal distribution)
- crosstabs コマンドとそのオプション

4.3. 無関連状態と期待度数

- 係数
- 期待度数・残差・連関係数
- クロス表とグラフの書きかた

5. 中間試験

6. 記述統計 (2): 平均値の比較

6.1. 平均と分散

- データの種類: 復習
- 順序尺度と間隔尺度の変換
- 平均値
- 分散と標準偏差
- 分布と外れ値

6.2. 平均値の層別比較

- 層別平均
- エフェクト・サイズ
- 相関比から分散分析へ
- 表とグラフの書きかた

7. 推測統計

7.1. 誤差の評価

- データの記述と誤差の評価
- 標本抽出の4段階モデル
- 無作為抽出
- 非標本誤差
- 標本誤差の統計的推測

7.2. 平均値の推定

- 平均値の点推定
- 区間推定と t 検定
- 平均値の差の区間推定
- エフェクトサイズ・相関比と区間推定

7.3. 統計的検定

- 区間推定の簡易表記としての有意水準
- 平均値の差の t 検定
- 連関係数の 2 検定
- 分散分析と F 検定
- 検定結果の表記

8. 期末レポート

2009.4.14

現代日本論演習

統計分析の基礎

東北大学文学部 2009 年度
田中 重人 (講師)

1

【目的】

統計分析の基礎的な手法の習得

SPSS の操作

クロス表分析

平均値の比較

推測統計の手法

2

【教科書】

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。

3

【成績評価】

- ・ 授業中の課題 (50%)
- ・ 中間試験 (20%)
- ・ 期末レポート (30%)

4

【関連する授業】

6 セメスタ

・ 現代日本論演習「質問紙法の基礎」
(水2)

・ 現代日本論演習「実践的統計分析法」
(木2) ...大学院と合同

5

受講登録フォーム記入

6

【コンピュータ実習室について】

入室に**学生証**が必要

土足・飲食・喫煙 **厳禁**

退出時は必要事項を紙に書く

(書けるところを書いてみよう)

ドアの開けかた

7

【コンピュータの起動と終了】

- ・ 本体とディスプレイの電源を ON
- ・ 表示されるお知らせの内容をよく読む
- ・ 終了するとき、ディスプレイの電源を切ることをわすれないように

8

【ファイルの保存場所】

授業でつかうファイルは、
授業開始時に マイドキュメント
フォルダにコピーして使う。
授業終了時に削除してかえること。

内蔵 Disk にデータは置けない

9

必要なデータは各自で
フロッピーかスティックメモリ
にコピーして持ち帰る

各自で購入しておくこと。

10

【SPSS】

データ解析用ソフトウェア

Windows での開発に

特に力を入れている

購入しやすい

11

【この授業で使用するデータ】

1995 年 SSM 調査 B 票の一部

cf. 『日本の階層システム』(全 6 巻)
東京大学出版会、2000 年。

SSM 調査については <http://www.sai.tohoku.ac.jp/coe/ssm/> 参照

12

2009.4.14

現代日本論演習 (田中重人)

受講登録フォーム

氏名：

学年：

学生番号：

所属 (文学部日本語教育学専修以外の場合)：

研究内容：

- ・ 自宅でパソコンを使えますか? **ある / ない**
- ・ SPSS を使った経験がありますか? **ある / ない**
- ・ コンピュータ・プログラムを作成したり、プログラミングの授業を受けたりしたことがありますか? **ある / ない**
 ある場合 **言語名** ()

数学的予備知識の調査 (成績評価には関係ありません)

(1) 「乱数」とは何か。簡単に説明せよ。

(2) 「必要十分条件」とは何か。簡単に説明せよ。

(3) 「偏差値」はどのような目的のために使われるか。またどうやって求めるか。簡単に説明せよ

(4) つぎの数式の値を求めよ。計算のプロセスがわかるように解答すること

$$\sum_{k=1}^{10} k =$$

受講者の興味と数学的知識の調査

別紙

コンピュータ実習室について

入室・退室

学生証が必要 (ない人は、教務係で臨時カードを借りること)。

土足・飲食・喫煙厳禁。

退出時には必要事項を紙に記入。

コンピュータの起動と終了

使いはじめるときは……

- コンピュータ本体の電源を入れる
- ディスプレイの電源を入れる (2-3秒押しつづけないと入らないので注意)
- 表示されるお知らせをひととおりよむこと
- キーボード右上の「NumLock」ランプがついているか確認

使い終わるときは……

- 「マイドキュメント」などに保存してある自分のファイルを削除
- 画面左下の「スタートメニュー」から「終了オプション」「電源を切る」を選択
- コンピュータ本体の電源が切れたことを確認
- ディスプレイの電源を切る
- フロッピーディスク、USBスティック・メモリなどをわすれないこと

ファイルの保存場所について

教室のコンピュータの内蔵ディスクには、個人のファイルを置いてはならない。授業中に必要なファイルは「マイドキュメント」フォルダに一時的に保存してよいが、授業が終わったら自分のフロッピーかスティック・メモリ等にコピーして、内蔵ディスクのほうのファイルは削除すること。

コンピュータ実習室で使えるリムーバブルメディアはつぎのふたつ。各自どちらかを購入しておくこと。

- フロッピーディスク (3.5インチ) ……「Windows フォーマット」のものが便利。安いがよく故障する。容量が小さい。
- フラッシュメモリ ……「USB2.0対応」のもの。値段は高いが容量が大きい。とりはずすときは画面右下の「ハードウェアの安全な取り外し」アイコンをクリックして、「USB大容量記憶装置」を停止させてから、メモリ本体を引き抜く。

模擬データ入力実習

SPSS について

参考書: 宮脇典彦・和田悟・阪井和男 (2000) 『SPSSによるデータ解析の基礎』培風館。

SPSS の起動

スタートメニューから「プログラム」「SPSS for Windows」「SPSS for Windows 12.0J」で起動する。(ここで何かエラーメッセージが出るかもしれないが、気にせず「続行」または「OK」する。)

「どのような作業を行いますか?」ときかれたら「データを入力」をチェックして「OK」。

データ入力

配布した架空の回答票をもとに、データを入力してみよう。

まず変数を定義

- 「データエディタ」ウインドウのいちばん下の「変数ビュー」タブに切り替える
- 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らんで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「マイドキュメント」内に保存してみる。(ほかのフォルダに保存してはならない。)

「マイドキュメント」を開いて、SPSS 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、コピーしておく必要はない。)

この方式はSPSS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ので、テキストファイルでデータを用意しておくのがふつうである。

1. データの配布
2. 標本抽出
3. SPSS のウインドウ構成
4. 変数値の再割り当て
5. 出力の読みかた・印刷

1

【データの配布】

1995 年 SSM 調査 B 票の一部
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
cf. (2000) 『日本の階層システム』(全 6 巻)
東京大学出版会。

2

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
流出しないように
変数ラベルは菅野剛
(日本大学) 氏による

3

毎回の授業で使うので、
忘れないこと (調査票も)
期末レポート提出時に返却

4

【標本抽出の 4 段階モデル】

ユニバース (universe)
母集団 (population)
計画標本 (designed sample)
有効標本 (valid sample / case)

5

伝統的な統計学では 4 段階に
わけずに、2 段階で考えるのが
ふつう：

母集団=universe + population
標本 = (designed/valid) sample

6

【無作為抽出】

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 (random sampling)
→ 「**等確率標本**」

7

つぎの条件が必要：
母集団の人口が既知
個体を網羅した「台帳」

個体によって抽出確率が違う場合も、事後的に調整して
等確率標本と同様の統計処理をおこなうことは可能
「台帳」が完備してない状況でも、工夫次第で
無作為抽出に近づけることができる

8

統計的な推測は、**等確率標本を前提とする**

実際の調査で理想的な標本抽出ができることはまずない。
また計画標本のなかから無効回答があるので、
無作為ではない誤差がかならず発生する。
この誤差は **統計的には処理できない**ので、個別に推測する

- ・ どの層を過剰に代表しているかを把握する
- ・ おなじ母集団を対象にした調査と比較する

9

【層化 2 段無作為抽出】

- ・ 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)

10

【データ・セット】

ケース × 変数
変数は変数名で管理
変数名以外に「ラベル」
無回答などの欠損値 (.)

11

【SPSS のウインドウ構成】

データ・エディタ
シンタックス・エディタ
出力ビューア

12

【メニューとシンタックス】

分析手法をえらぶ
必要なオプションを指定
「貼り付け」をクリック
シンタックスの必要部分を選
択して実行 (▶)

13

【出力ビューア】

左側に目次、右側に出力内容
エラー表示もここに出る

【印刷】

左側の目次で選択
出力先の切り替え
印刷前にプレビュー
電源の入れかた
ジョブの確認・取り消し
タイトル印刷 (2 面, 4 面, ...)

14

【変数値の再割り当て】

データエディタのメニューバーで
「変換」 「値の再割り当て」
「他の変数へ」
変換先変数の名前をつける

15

「今までの値と新しい値」
値の組を指定したら「続行」
シンタックスを貼付けて実行
新変数の度数分布を確認
問題がなければデータセット
を保存

16

1. 尺度水準
2. 度数分布表

1

【尺度水準】

比率尺度 (ratio scale)
間隔尺度 (interval —)
順序尺度 (ordinal —)
名義尺度 (nominal —)
(質的変数とも)

(教科書 p. 8)

2

【尺度の変換】

上位の尺度のほうが
あつかえる演算が豊富
上位の尺度は下位の尺度の特
徴を兼ね備えている

分析手法の選択幅がひろい

3

私たちが測定するものはいいてい
順序尺度以下である

上位の尺度への変換には
一定の理論的根拠が必要

4

【実習】

SSM 調査の調査票中で、比率尺度
とみなせるものはどれか

5

【度数分布表】

Frequencies コマンド

「分析」
「記述統計」
「度数分布表」

6

出力：

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)

7

【累積%とパーセンタイル】

順序尺度以上の場合のみ意味を持つ
Percentile(= %点)
中央値 (median) = 50%点
「割り切れてしまう」場合は中点をとる
(教科書 p. 43)

同じ値が並ぶ場合は多少の操作が必要
(森敏昭・吉田寿夫(編)(1990)『心理学のための
データ解析テクニカルブック』北大路書房. p. 15)

8

【実習】

世帯収入 (q44_3) について、度数分布表を出
力し、中央値、25%点、75%点を求めよ

9

- 1. グラフの利用
- 2. 棒グラフとヒストグラム

1

【グラフの利用】

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

2

【棒グラフとヒストグラム】

棒グラフ.....棒同士の間空白をあける。高さ(長さ)をよむ。
histogram (柱グラフ).....柱の間隔をあけない。面積をよむ。

縦軸は度数または%

3

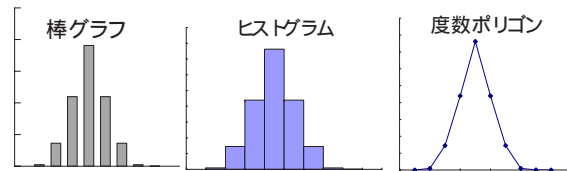
連続量を階級分けした場合
ヒストグラム

それ以外の場合 (離散量 /
名義尺度) 棒グラフ

度数多角形 (polygon) は複数の変数の分布を比較するときに便利。

(教科書 p. 32-36)

4



SPSS では histogram が書きにくい。

recode で整形した上で度数分布表のメニューで「図表...」指定。棒グラフを書く
グラフ インタラクティブ ヒストグラム
では等間隔の区間に分割してくれる

5

Excel を使う場合：

recode で整形した上で度数分布表を出力
表を Excel にコピーする

棒グラフを作成

グラフの棒の上で左クリック

「データ系列の書式設定」

「オプション」

「棒の間隔」を0にする

6

【実習】

適当な変数について棒グラフまたはヒストグラムを作成

7

【キーワード】

行 (row) 列 (column) セル (cell)

周辺度数 (marginal frequency)

行% (row percent) 列% (column percent)

1

【度数分布表の比較】

データエディタのメニューで
「データ」「ファイルの分割」
「グループの比較」

度数分布表を出力

2

「データ」「ファイルの分割」
「すべてのケースを分析」
でもとにもどしておく

3

【クロス表の基本型】

質的変数 (名義尺度) 同士の関連
についての基本的な分析法

	1	2	3	合計
行 1	a	b	c	a+b+c
2	d	e	f	d+e+f
3	g	h	i	g+h+i
合計	a+d+g	b+e+h	c+f+i	N

列

周辺度数

4

5

【Crosstabs コマンド】

性別 × 「性別による不公平」
のクロス表を書いてみよう

「分析」「記述統計」「クロス集計表」

6

【行%と列%】

「クロス集計表」メニューで「セル」にパー
センテージ (行・列) を追加

行%, 列%のつかいわけは

説明 被説明の関係に対応

行 列の説明をすることが多い

周辺度数の%とも比較する

7

【グラフを書いてみる】

クロス表は帯 (積み上げ棒)

グラフで表現することが多い

SPSS ではうまくかけない。コピーして
Excel に貼付けてグラフを書くのがよい

度数にも注意

8

【課題】

性別 × 適当な変数でクロス表作成、
%からわかることをコメントする。
グラフも書いて印刷して提出

9

第7回「φ係数」

1. 自由度 (degree of freedom)
2. クロス表分析のふたつの系列
3. 2×2 クロス表の性質
4. φ係数 (phi coefficient)

1

【自由度】

2×2 クロス表では、周辺度数が所与なら、
1つのセル度数が決まればほかも決まる

α	β		合計
	1	2	
1	a	g-a	g
2	i-a	h-i+a	h
合計	i	j	N

2

3×3 クロス表：セル度数が4つ決まれば…

α	β			合計
	1	2	3	
1				f
2				g
3				h
合計	i	j	m	N

k×l クロス表の自由度 (degree of freedom)

$$d.f. = (k-1)(l-1)$$

3

【クロス表分析の2つの系列】

- 「%の差」系 (期待度数との差)
= 連関係数
- オッズ比系 (乗法モデル)
= 対数線形分析、ロジット分析

この授業で取り上げるのは前者だけ

4

【2×2 クロス表の性質】

以下、つぎの記号法を使う

α	β		合計
	1	2	
1	a	c	g
2	b	d	h
合計	i	j	N

5

(1) 行%は1列について比較すればよい:

$$\frac{a}{g} - \frac{b}{h} = \frac{d}{h} - \frac{c}{g}$$

(2) 行%の差がゼロなら列%の差もゼロ

(3) 行%の差が100なら列%の差も100

(4) $g=i$ or $g=i$ なら行%の差と列%の差は同じ:

$$\frac{a}{g} - \frac{b}{h} = \frac{a}{i} - \frac{c}{j}$$

6

(5) これら以外の場合、行%の差と列%の差はちがう値になる

(例1) 行%の差 = 8%

60%	40%	100%
52%	48%	100%

(例2) 行・列とも%に差なし

52	48	100
52.0%	48.0%	100.0%
66.7%	66.7%	
26	24	50
52.0%	48.0%	100.0%
33.3%	33.3%	
78	72	150
52.0%	48.0%	100.0%

(例3) 行・列とも10%の差

70	30	100
70.0%	30.0%	100.0%
70.0%	60.0%	
30	20	50
60.0%	40.0%	100.0%
30.0%	40.0%	
100	50	150
52.0%	48.0%	100.0%

8

【φ係数】

2×2 クロス表の「連関」の尺度

$$\phi = \frac{ad - bc}{\sqrt{ghij}}$$

この係数の意味は?

(分子だけ取り出して考えてみよう)

9

【キーワード】

連関 (association), 独立 (independence),
期待度数 (expected frequency),
クラメールの連関係数 (Cramer's V)

1

【φ係数の性質】

- φ = 交差積の差 / √(周辺度数の積)
- φ = 相関係数の特殊ケース
(→ VI セメスタ授業)
- |φ| = 行%差と列%差の中間の値
(教科書 p. 103 表 4-1 について計算してみよう)

2

4. φ² = 標準残差の2乗の総計 / N
(→ 2×2 以上のクロス表に拡張できる)

3

【期待度数とφ係数】

※記号法は前回と同じ

独立 (無関連): a/b = c/d

期待度数 (expected frequency)
周辺度数を固定しておいて独立なクロス表を作ったとき、各セルに入る度数:

$$\frac{gi/N}{hi/N} \quad \frac{gj/N}{hj/N}$$

4

各セルの期待度数は?

		100	
		100.0%	
		50	
		100.0%	
78	72	150	
52.0%	48.0%	100.0%	

5

- ★ 期待度数はたいてい小数になる
- ★ 期待度数について行%と列%を計算すると、周辺度数の%とおなじになる

観測度数 各セルに入る実際の度数
残差 (residual) 観測度数と期待度数の差
標準残差 (standardized ---) 残差/√期待度数

ex. $A = \frac{a - gi/N}{\sqrt{gi/N}}$

6

観測度数が下記の場合、
各セルの残差と標準残差は?

40	60	100
		100.0%
38	12	50
		100.0%
78	72	150
52.0%	48.0%	100.0%

7

χ² (chi-square) 標準残差の平方和

各セルに入る標準残差を A, B, C, D とする

$$\chi^2 = A^2 + B^2 + C^2 + D^2 = N \left(\frac{a^2}{gi} + \frac{b^2}{hi} + \frac{c^2}{gj} + \frac{d^2}{hj} - 1 \right)$$

χ² を人数で割った値が φ の2乗 に等しい

$$\phi^2 = \frac{\chi^2}{N} \quad \text{すなわち} \quad |\phi| = \sqrt{\frac{\chi^2}{N}}$$

8

【クラメールの連関係数 V】

k×l 表へのφ係数の拡張 (教科書 p. 114-117)

- ★ k と l のうち小さいほうを m とする
- ★ 2×2 表と同様に期待度数・残差を求める
- ★ χ² を求める
- ★ χ² を N と (m-1) で割って平方根をとる

$$V = \sqrt{\frac{\chi^2}{N(m-1)}}$$

9

【Vの性質】

- ★ 行・列変数が独立のとき V = 0
- ★ 関連が強くなると大きくなる
- ★ 最大値は 1

10

【SPSSで実習】

クロス表のオプションを指定:

- 「セル」… 度数(観測/期待)
残差(標準化なし/標準化)
- 「統計」… カイ2乗
ファイと Cramer の V

11

【課題】

教科書 p. 103 表 4-1 について、
期待度数、残差、標準残差、χ²、Cramer の連関係数 V
を計算。

12

【予告】

再来週 (6/30) は中間試験

- ・ 何でも持ち込み可
- ・ 出題範囲は、6/23 授業まで

13

第9回「平均値と標準偏差」

1. 尺度水準と分析法
2. 代表値と散布度
3. 平均値と標準偏差
4. SPSS のコマンド
5. 平均値の層別比較

1

【尺度水準と分析法】

名義×名義 クロス表

名義×間隔 平均値の比較

2

【代表値と散布度】

中央値 (median) - 四分位偏差 (Q)
(順序尺度以上)

平均値 (mean) - 標準偏差 (SD)
(間隔尺度以上)

(教科書 p.42-51)

3

【平均値】

総和をデータ数で割ったもの

【標準偏差】

平均値からの偏差の2乗値の平均が「分散」
分散の平方根が「標準偏差」

平均値と標準偏差はセットで使う

4

次のデータの平均とSDは?

{0, 1, 4, 5, 7}

5

値	偏差	偏差 ²
0		
1		
4		
5		
7		

平方和 =

分散 =

SD =

6

【SPSS のコマンド】

「記述統計」 「度数分布表」

「統計」オプションで
「平均値」と「標準偏差」をチェック

「記述統計」 「記述統計」でもよい

7

【平均値の層別比較】

ふたつの層の間の平均値の比較

平均値の差をもとめる
(層別平均)

標準偏差を基準にして差を評価
(effect size) 次回

8

【SPSS のコマンド】

「平均の比較」 「グループの平均」

従属変数 = 平均値を求める変数
(間隔尺度)

独立変数 = 層を指定する変数
(名義尺度)

9

1. 平均値使用時の注意事項
2. エフェクト・サイズ

1

【平均値を使うときの注意事項】

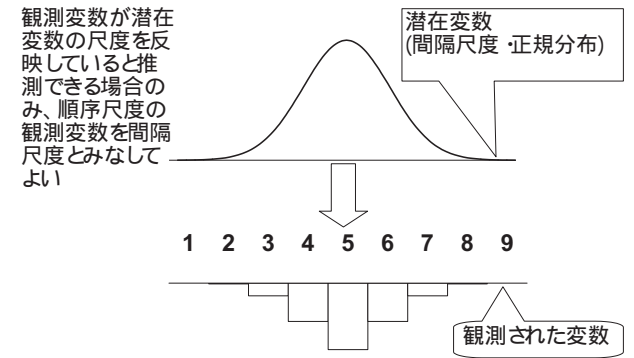
順序尺度の平均値をとっていいのは

- ・ 潜在的には間隔尺度のはず
- ・ 測定のポイントが一定間隔

という 2 条件をともに満たす場合

2 値の変数は間隔尺度とみなせるが、若干の注意が必要。

2



3

具体的には

4 点以上の尺度

正規分布に近似 (教科書 p. 53-59) :

- ・ 単峰性
- ・ 左右対称性 (歪度)
- ・ 中央への集中度 (尖度)

ヒストグラムを描いて検討するとよい。

正規分布との乖離度を統計的に検討する手法もある

4

歪度・尖度は「度数分布表」の「統計」オプションで指定できる

正規分布のとき 0、絶対値が大きくなるほど、正規分布から外れる

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)

順位に変換したり中央値を使って分析

5

平均値ははずれ値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

左右対称でないデータ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

6

【エフェクト・サイズ】

$$ES = \text{平均値の差} / \text{標準偏差}$$

正式には層別 SD の重みつき平均のような数値 (併合 SD) をつかう (教科書 p. 137)

7

【例】

性別による生活全般満足度の違い

	平均	SD	(人数)
男性	2.62	1.02	(114)
女性	2.24	0.91	(136)
合計	2.41	0.98	(250)

平均の差 =
併合 SD
ES =

ES は SPSS では計算してくれない

8

【期末レポート】

期限: 8/10 (月) 17:00

提出先: 日本語教育学研究室 (文法合同棟 2F)
205 室の田中のレターケース

内容: クロス表と平均値の比較について適当な分析をして結果を解釈する。推測統計の結果をかならずふくめる。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。

備考: 後期の授業を受講しない者は、SSM データのディスクをレポートと一緒に提出。データのコピーをすべて消去すること。

9

1. 相関比

2. 分散分析 (ANOVA)

3. エフェクト・サイズと相関比

1

【ES の特徴と問題点】

各層の人数を考慮せず平均値だけ比較

➔ 大きさがちがう場合は？

2 層間の比較だけ

➔ 3 つ以上の層を比較したい場合は？

2

【相関比】

各層の個体が全員その層の平均値を持つ
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」。(イータ) であらわす

相関比の 2 乗 ² を

「決定係数」「分散説明率」などという

² を「相関比」ということもある

3

SPSS では

「オプション」の「第 1 層の統計」で
「分散分析表とイータ」をチェック

は 0~1 の範囲の値をとり、
独立変数の影響力をあらわす

ES は最小値 0、最大値

4

3 層以上で平均値を比べる場合にも
相関比が使える。

このように、層別平均値をあてはめて仮想
分散を求める分析法を「**分散分析**」
(ANOVA: ANalysis Of VAriance) という。

5

【ES と の関係】

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

特に、2 層の大きさが同じ ($n_1 = n_2$) なら、

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

層の大きさがちがえば、ES はこれより大きくなる

6

このように ES と は互いに変換できる。

両方示すのは冗長

7

【注意事項】

層別の平均値を分析する場合、
各層の人数は一定以上必要

(最低 20 人?)

カテゴリ統合が必要になることがある

8

【課題】

(1) 適当な変数の平均値について、
男女別の平均値の差と ES、 を求める
(併合 SD のかわりに層別 SD の単純平均を使ってよい)。
表に ES と を書き込んで提出。

(2) 平均と SD の表 (前回資料) から の近似
値を求める方法を考えて書いてくる。

9

現代日本論演習 (田中 重人)
2009.7.14 課題

氏名:
学年:
所属:
学生番号:

次の3つの表の網掛け部分を埋めよ

全体についての平均と標準偏差			
	平均	偏差	偏差 ²
1	3.00	-2.00	4.00
1	3.00	-2.00	4.00
2	3.00	-1.00	1.00
2	3.00	-1.00	1.00
3	3.00	0.00	0.00
5	3.00	2.00	4.00
4	3.00	1.00	1.00
5	3.00	2.00	4.00
4	3.00	1.00	1.00
3	3.00	0.00	0.00
合計	30	30.00	
平均	3.00	3.00	
SD=			

層別平均を当てはめた仮想データセットの平均と標準偏差			
層別平均	全体平均	偏差	偏差 ²
1.50	3.00		
1.50	3.00		
1.50	3.00		
1.50	3.00		
4.00	3.00		
4.00	3.00		
4.00	3.00		
4.00	3.00		
4.00	3.00		
4.00	3.00		
4.00	3.00		
合計	30.00	30.00	
平均	3.00	3.00	
SD=			

層別の平均と標準偏差			
層別平均	偏差	偏差 ²	
1	1.50		
1	1.50		
2	1.50		
2	1.50		
合計	6	6.00	
平均	1.50	1.50	
SD=			
3	4.00		
5	4.00		
4	4.00		
5	4.00		
4	4.00		
4	4.00		
3	4.00		
合計	24	24.00	
平均	4.00	4.00	
SD=			

=
^2 =

現代日本論演習 (田中 重人)
2009.7.14 課題

氏名:
学年:
所属:
学生番号:

次の3つの表の網掛け部分を埋めよ

全体についての平均と標準偏差			
	平均	偏差	偏差 ²
1	3.00	-2.00	4.00
1	3.00	-2.00	4.00
2	3.00	-1.00	1.00
2	3.00	-1.00	1.00
3	3.00	0.00	0.00
5	3.00	2.00	4.00
4	3.00	1.00	1.00
5	3.00	2.00	4.00
4	3.00	1.00	1.00
3	3.00	0.00	0.00
合計	30	30.00	0.00
平均	3.00	3.00	2.00
SD= 1.41			

層別平均を当てはめた仮想データセットの平均と標準偏差			
層別平均	全体平均	偏差	偏差 ²
1.50	3.00	-1.50	2.25
1.50	3.00	-1.50	2.25
1.50	3.00	-1.50	2.25
1.50	3.00	-1.50	2.25
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
4.00	3.00	1.00	1.00
合計	30.00	30.00	0.00
平均	3.00	3.00	15.00
SD= 1.22			

層別の平均と標準偏差			
層別平均	偏差	偏差 ²	
1	1.50	-0.50	0.25
1	1.50	-0.50	0.25
2	1.50	0.50	0.25
2	1.50	0.50	0.25
合計	6	6.00	0.00
平均	1.50	1.50	0.25
SD= 0.50			

3	4.00	-1.00	1.00
5	4.00	1.00	1.00
4	4.00	0.00	0.00
5	4.00	1.00	1.00
4	4.00	0.00	0.00
3	4.00	-1.00	1.00
合計	24	24.00	0.00
平均	4.00	4.00	0.67
SD= 0.82			
併合SD = 0.71			

= 0.866
^2 = 0.750

全体SD² = 2.00
仮想SD² = 1.50
併合SD² = 0.50

平均

第 13 回「推測統計の基礎」(2009.7.21)

1. 記述統計と推測統計
2. 無作為抽出
3. 点推定と区間推定
4. 比率の区間推定

1

【記述統計と推測統計】

記述統計 (descriptive statistics)

= データ (ケース) の特徴を
数値や図表にまとめる

推測統計 (inferential statistics)

= 確率的な誤差を考慮して、
母集団の特徴を推測する

(教科書 pp. 3-5)

2

【無作為抽出】

random sampling

母集団から計画標本を選ぶ際に、

すべての個体の抽出確率が等しくなる

ように抽出する

→ 「等確率標本」(probability sample)

3

袋のなかに色つきの玉が 60 個入っている:

赤色: 30 個

青色: 20 個

黄色: 10 個

玉を n 個取り出したとき、その色は.....?

全世界から n 人を無作為抽出したとき、
そのなかに 人は何%ふくまれるか?

4

【区間推定】

interval estimation

「答えは たぶん この範囲内にある」

信頼率 (confidence level) を適当に設定して

信頼区間 (confidence interval) を求める

5

全世界から 400 人を無作為抽出:

うどん が好き: 240 人

そば が好き: 160 人

うどんが好きな人の比率は?

$$0.6 \pm 1.96 \times (0.6 \times 0.4 / 400)$$

答: _____ ~ _____ % (95%信頼区間)

6

【比率の区間推定】

標本の規模がじゅうぶん大きく ($n > 30$),

比率があまり偏っていない ($0.1 < m < 0.9$) とき、

95%信頼区間は

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

標準誤差
(standard error)

7

【平均値の区間推定】

母集団における平均値についても、同様の計算ができる。
ただし、正規分布を仮定:

$$\frac{m}{\text{標本平均}} \pm 1.96 \times \frac{\text{SD}}{\sqrt{n}} \text{ (標準誤差)}$$

↓
t 臨界値

「t 臨界値」は n によって変化するが、 $n > 200$ で 1.96 に収束 (教科書 p. 281)。

8

【SPSS コマンド】

「分析」 「記述統計」 「探索的」

「従属変数」を指定
パネル左下の「統計」だけをチェック

信頼率を変更するには「統計」を選択
「因子」を指定すると層別に分析できる

9

1. 平均値の差の推定
2. 区間推定と統計的検定

1

【平均値の差の推定】

2 層間の **平均値の差** についても
 平均値そのものと同様の区間推定ができる：
 このとき 95%信頼区間は

$$d \pm t_{\text{臨界値}} \times \text{併合SD} \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

標準誤差

ただし n_1, n_2 はそれぞれの層の人数
 t 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

2

【SPSS のコマンド】

「平均値の比較」 「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。
 連続量を一定の値で切ることできる

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

3

【統計的検定】

Statistical test

統計的検定 = 特定の値を設定して、その値が
 信頼区間に含まれているかどうかを判定する
 0 に設定するのがふつう

4

【統計的検定用語】

帰無仮説 (null hypothesis):

母集団における統計量が
 この「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

危険率 (critical level): 1 - 信頼率

5

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 0 をふくまない
 = 少なくとも 95%の確率で、
 母集団において平均値の差がある
 といえる

6

「5%水準で非有意」とは.....
 95%信頼区間が 0 をふくむ
 = 母集団においては平均値の差がない
 という確率が 5%以上ある

7

【有意確率とは】

信頼区間をひろげていくと、
 どこかでゼロ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。

8

分析の際は、

- ・ 前もって危険率を設定しておく
(通常は 5%または 1%)
- ・ 有意確率はその値を
下回っているかどうか判別する

例:

有意確率が 0.007
 有意確率が 0.023
 有意確率が 0.088

9

【区間推定と統計的検定】

区間推定と統計的検定の方法の間に本質的なちがいはない

慣習的に統計的検定を使うことが多い(分野によってちがう)

統計量によっては、区間推定はすごくむずかしい場合がある

10

【むずかしい区間推定】

係数 「Fisher の z' 変換」をおこない標準正規分布を利用 (相関係数と同じ) 森・吉田 (1990, p. 225)

連関係数 V 非心 χ^2 分布を利用

相関比 非心 F 分布を利用

11

【平均値の差の t 検定】

コマンドの指定は区間推定とおなじ。出力の「有意確率 (両側)」を見る

2 層の間の差の検定にしか使えない
「母集団では正規分布」を前提とする
2 層の間で分散が等しいことが前提

12

【クロス表の独立性の検定】

V または χ^2 の信頼区間にゼロ (=独立の状態) がふくまれるかを判別する。

「クロス集計表」の「統計」で「カイ 2 乗」を指定。

出力の「Pearson」の列の右端が有意確率

各セルの期待度数が 5 以上であることを前提とする

13

【分散分析と F 検定】

「平均値の比較」「グループの平均」オプション「分散分析表とイータ」を指定
出力「分散分析表」の右端「有意確率」

3 層以上の場合に使う。
の信頼区間を使って判断するのと同じである。
2 層の場合にも使えるが、 t 検定と同じ結果になる必要とする前提も t 検定と同様

14

【表の書きかた】

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
 $p < 0.05$ 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)

15

【文献】

森敏明・吉田寿夫 (1990) 『心理学のためのデータ解析テクニカルブック』北大路書房。

16

2009.7.28 現代日本論演習 (田中重人)

授業資料

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)

=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)

= 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)

= 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)

= 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)

= 0.086^{ns}. ns: 5%水準で非有意.

無回答 = 7.