

社会調査法 第5回

1. 基本統計量

次の表に、このクラスの学生のデータを書いてみましょう。

身長は平均どのくらいだと思いますか？男性と女性でどのくらい変わりますでしょうか？

学生	性別	身長 (cm)	今アルバイトをやっているか (SA)	アルバイトをしている場合、週の日数	アルバイトをしている場合、アルバイトの種類 (MA 複数回答)
1	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他
2	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他
3	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他
4	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他
5	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他
6	1 女性 2 男性 3 その他	cm	1. いつもやる 2. 時々やる 3. やらない	日	1 サービス業 (飲食、販売など) 2 製造業 (弁当製造、組み立てなど) 3 教育・福祉 (家庭教師や老人ホーム) 4 その他

2. 重要な用語

(1) 標本誤差

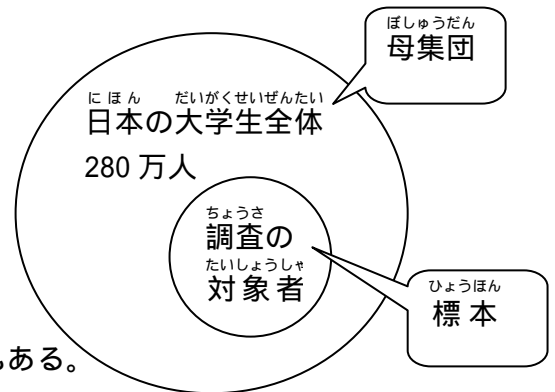
$$\text{標本誤差} = 1.96 \sqrt{\frac{N-n}{N-1} \times \frac{P(100-P)}{n}}$$

母集団が大きいときは右でもいい(あまり変わらないから) 標本誤差 = $2 \sqrt{\frac{P(100-P)}{n}}$

n = 調査の対象者数 (標本数)

P = 回答率

N = 母集団の人数



(シグマ 母集団がはっきりしている場合)を

標準偏差と、SE または S (標準誤差) という場合もある。

例: もし無作為に選んだ 80 人の大学生を対象者に、「アルバイトをしているか」聞いて、アルバイトをしている人が 60 人だった場合 (75% の時) 標本誤差 () は

$$\sigma = 2 \sqrt{\frac{P(100-P)}{n}} \quad \text{つまり} \quad \sigma = 2 \sqrt{\frac{75(100-75)}{80}} = \text{つまり } \boxed{5} \text{ です。}$$

これはアルバイトをしている人の割合が、母集団 (日本の大学生全体 280 万人) 70% から 80% の間 (75% ± 5%) ということです。

±5% では、標本誤差が大きすぎるので、実際にはもっとたくさんの人を調査対象者にした方がいいでしょう。

また、正確な結果を得るには、対象者の全員に回答してもらうことが必要です。しかし現実の調査では、全員回答してもらうのは、非常に難しく、行政が行う調査でも回答率は 50% 位にしかならないのが現状です。

(2) 平均

身長や金額などの数値データの時のみ、集計できます。

例: もしこのクラスの学生の身長が

156cm	182cm	161cm	172cm	161cm	159cm	173cm	169cm
-------	-------	-------	-------	-------	-------	-------	-------

のとき、

$$(156\text{cm} + 182\text{cm} + 161\text{cm} + 172\text{cm} + 161\text{cm} + 159\text{cm} + 173\text{cm} + 169\text{cm}) \div 8 \text{ 人で } \text{平均は } 166.6\text{cm}$$

(3) 分散 ^{ぶんさん}

データのばらつきを表します。(個々のデータ - 平均)の2乗の総和(全部足した数) ^{ここ} ^{へいきん} ^{にじょう} ^{そうわ} ^{ぜんぶ} ^た ^{かず}
 分散が大きければ、ばらつきが大きい。(元のデータの絶対値が大きければ大きくなる) ^{ぶんさん} ^{おお} ^{おお} ^{ぜったいち} ^{おお}

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

は「総和」と読む ^{そうわ}

\bar{X} は「エクスペー」または「平均」と読む ^{へいきん}

例	1	2	3	4	5	6	7	8
身長	156cm	182cm	161cm	172cm	161cm	159cm	173cm	169cm
A 身長 - 平均 (166cm)	-10cm	16cm	-5cm	6cm	-5cm	-7cm	7cm	3cm
B A × A (二乗)	100cm	256cm	25cm	36cm	25cm	49cm	49cm	9cm
Bの総和	549cm							
B ÷ 8 (人数)	69	これが分散						

(4) 標準偏差 (SD standard deviation) ^{ひょうじゆんへんさ}

偏差(個々のデータ - 平均値)の2乗の和(足し算したもの)を、データの数で割り、平方根(ルートのこと)をとる。 ^{へんさ} ^{へいきんち} ^{にじょう} ^わ ^た ^{ざん} ^{かず} ^わ ^{へいほうこん}

データのばらつき具合がわかる。(分散と違って、元のデータの単位となっている) ^{ぐあい} ^{ぶんさん} ^{ちが} ^{もと} ^{たんい} ^{たとえ}
 ばらつきが上下8.3cmとなる。分散では69) ^{しんちよう} ^{じようげ}

$$SD = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}}$$

なぜ、わざわざ二乗して、ルートをとるのか? ^{じじょう}
 平均より大きい+のデータ、平均より小さいマ ^{へいきん} ^{へいきん} ^{ちい}
 イナスのデータの、符号を無視するため ^{ふごう} ^{むし}

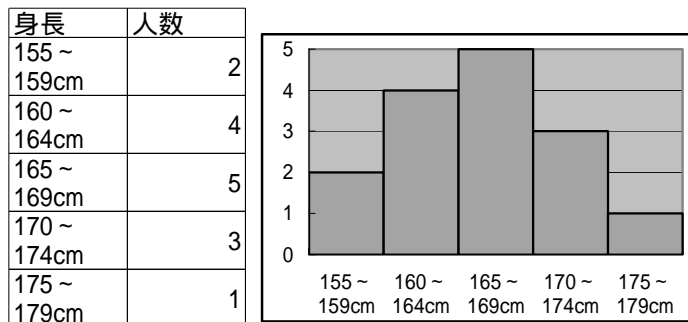
例 $\sqrt{69} = 8.3$ ^{れい} ^{だいたい} 8.3cmのばらつきがあるということ

	1	2	3	4	5	6	7	8	けい 計
身長 ^{しんちよう}	156cm	182cm	161cm	172cm	161cm	159cm	173cm	169cm	
A 身長 - 平均 (166cm) ^{しんちよう} ^{へいきん}	-10cm	16cm	-5cm	6cm	-5cm	-7cm	7cm	3cm	
B A × A (二乗) ^{じじょう}	100cm	256cm	25cm	36cm	25cm	49cm	49cm	9cm	
C Bの総和 ^{そうわ}									549cm
D 分散 C ÷ 8人 ^{ぶんさん} ^{にん}									69cm
E 標準偏差 D ^{ひょうじゆんへんさ}									8.3cm

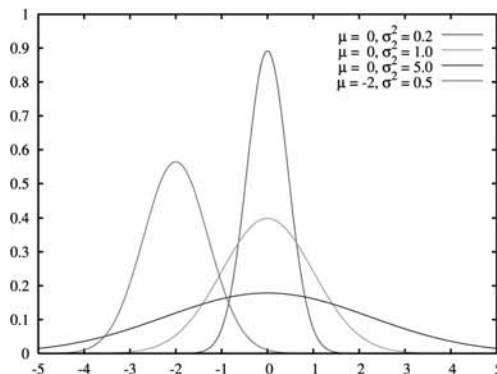
(5) さいだい値 さいしょう値 ちゅうおう値 さいひんち

1	2	3	4	5	6	7	8
156cm	159cm	161cm	161cm	169cm	172cm	173cm	182cm
さいしょう値		ちゅうおう値	ちゅうおう値				さいだい値
		さいひんち	さいひんち				
		最頻値	最頻値				

(6) どすうぶんぷひょう (ヒストグラム)



(5) せいきぶんぷ



Wikipedia より

(7) しゃくど

ひりつしゃくど ねんれい しゅうにゅう
比率尺度 年齢、収入など
 * 間隔尺度との違いは、絶対的規準「0」を持つこと

かんかくしゃくど しすう おんど けいさん
間隔尺度 指数、温度など * 計算できる 0 がない

れい ふかいしすう きおん しつど きおん
 例 不快指数 (= 0.81 × 気温 + 0.01 × 湿度 × (0.99 × 気温 - 14.3) + 46.3)

~ 55	55 ~ 60	60 ~ 65	65 ~ 70	70 ~ 75	75 ~ 80	80 ~ 85	85 ~
さむ 寒い	はださむ 肌寒い	なに かん 何も感じ ない	こちよ 快い	あつ 暑くない	ややあつ やや暑い	あつくてあせ 暑くて汗 で が出る	あつ 暑くてた まらない

じゅんじょしゃくど
順序尺度

とてもよい - よい - ふうふう - わるい - とてもわるい
など *計算できない

アルバイトをしているか
1 いつもやる 2 時々やる 3 やらない

ただし、わかりやすさや比較のために、次のように点数化して計算する場合もある

とてもよい - よい - ふうふう - わるい - とてもわるい
10点 5点 3点 2点 0点

めいぎ
名義尺度

性別など *計算できない

性別
1 男性 2 女性

名義尺度で、男性を1、女性を2などと番号を振って計算する場合「ダミー（模擬）変数」と呼ぶことがある。

(8) 変数の種類

りさんへんすう れんぞくへんすう
離散変数と連続変数

例：

性別が離散変数 男性、女性、その他 しかない

身長が連続変数 151、151.1 151.2 など小数点以下があり連続している

どくりつへんすう じゅうぞくへんすう
独立変数と従属変数

例：

学年によって、身長の高さが違う場合

学年が（身長などを決める）独立変数

身長が（学年によって決定する）従属変数

せつめいへんすう ひせつめいへんすう
説明変数と被説明変数

例：

性別で身長の平均が違うとき

性別が（身長を説明する）説明変数

身長が（性別によって説明される）被説明変数

	身長の平均
女性	161cm
男性	172cm

3. 相関関係とは

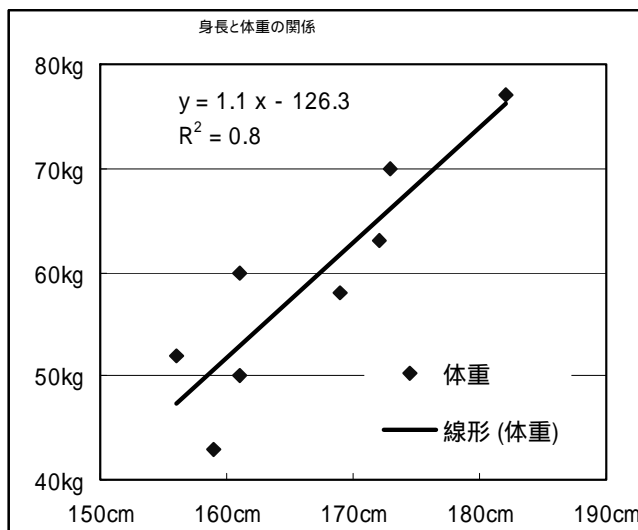
(1) 相関とは

ある変数とある変数の間に、相関関係があること。

例 身長と体重

	1	2	3	4	5	6	7	8
身長	156cm	159cm	161cm	161cm	169cm	172cm	173cm	182cm
体重	52kg	43kg	60kg	50kg	58kg	63kg	70kg	77kg

相関係数 0.89 + の相関(高ければ重い) 0.6以上は相関がある



(2) 疑似相関

相関があるように見えるが、実は他の要因と関係していること。

例：

誤り ジュースがたくさん売れると、水におぼれる人が増える？

正しくは 気温が上がると、ジュースが売れる、
水遊びをする人も増え、おぼれる人が増える

誤り？ 朝ご飯を食べると、成績がよくなる？

本当は？ 家庭環境で規則正しい生活ができる人は、朝ご飯を食べる人が多く、
成績もいいのかもしれない。