

心理・教育測定法 基礎編A

2017年度 心理：教育測定基礎編 A・B 心理：教育測定演習 A・B

心理・教育測定法(基礎編) A・B 木曜日 7-8限:W607			人間行動システム統計演習A・B 木曜日 9-10限 W607		
担当 前川			担当 前川 栗山		
4月 6日	オリエンテーション	前川	1	オリエンテーション	
4月 13日	統計量の計算・相関と回帰		2	データフレームの操作	
4月 20日	統計的推測の基礎		3	ヒストグラム、クロス集計	
4月 27日	一つの母平均に関する推測		4	相関と回帰	
5月 11日	二つの母平均に関する推測		5	平均値に関する推測	
5月 18日	複数の母平均に関する推測		6	分散分析1	
5月 25日	分散分析法と回帰分析		7	分散分析2	
6月 1日	まとめ		8	まとめ	
6月 15日	イントロ		1	イントロ	
6月 22日	多次元尺度法による個体の分類		2	多次元尺度法	
6月 29日	クラスター分析による個体の分類		3	クラスター分析	栗山
7月 6日	回帰分析		4	回帰分析	
7月 13日	因子分析		5	因子分析	
7月 21日	共分散分析と構造方程式モデル	栗山	6	共分散分析と構造方程式モデル	栗山
7月 27日	その他の方法		7	その他の方法	
8月 4日	まとめ		8	まとめ	

3

前川研究室 (Shin-ichi Mayekawa)

東京工業大学 大学院社会理工学研究科 人間行動システム専攻 前川研究室です。

<http://www.ms.hum.titech.ac.jp/>

の論文

uts

編 A B

統計的方法

- データの記述
 - 度数分布他
 - データの位置を示す指標 (点推定値)
 - データの散らばり度合いを表す指標
- 統計的推測
 - 現象のモデルとしての確率分布
 - 母集団からのサンプルとしてのデータ 標本分布
 - 母集団の位置に関する推測 統計的推定・検定
 - 複数母集団の比較
- 複数変数間の関係 → 基礎編 B

大事なこと：二つの世界

5

- 見える世界
 - データ、標本 (sample)
 - 標本統計量 (statistic)
- 見えない世界
 - 母集団、母集団分布
 - 母数 (parameta) 未知の定数
- 統計学の目的は
標本（可視）から母集団（不可視）を推測

データ

- 標本サイズ、サンプルサイズ、 N, n
 - 観測対象である個体の数、
number of observations, # of records
- 変数 variable
 - 観測した内容のこと
- 変数の数値的性質
 - 尺度水準
 - 連続か、離散か
 - 計測値か、計数值か

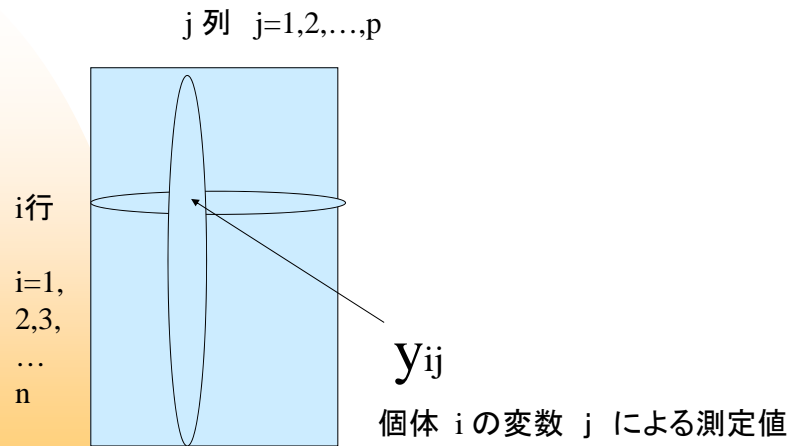
成績データ

ID	Gen	T1	Sco	GPA	ID	Gen	T1	Sco	GPA	ID	Gen	T1	Sco	GPA
1	1	165	60	3	18	1	114	56	3	35	1	155	46	3
2	1	127	62	3	19	0	140	45	4	36	1	85	66	3
3	1	146	50	4	20	1	91	70	4	37	0	221	86	4
4	1	109	47	2	21	1	93	60	3	38	1	116	45	3
5	0	72	19	3	22	1	159	44	3	39	0	118	51	4
6	0	70	40	3	23	0	141	28	2	40	0	110	40	4
7	0	88	37	3	24	1	131	70	3	41	1	123	52	3
8	0	33	37	2	25	1	143	67	4	42	1	146	46	5
9	1	196	71	5	26	0	224	90	4	43	0	89	29	2
10	1	142	72	3	27	0	126	77	3	44	1	144	69	5
11	0	170	61	4	28	0	202	63	5	45	0	91	46	3
12	1	217	77	5	29	1	93	49	2	46	0	55	0	1
13	0	106	59	2	30	0	54	9	2	47	1	52	53	3
14	1	128	42	3	31	1	108	34	2	48	1	157	57	3
15	1	32	14	1	32	1	92	32	3	49	0	93	34	3
16	0	152	57	2	33	0	133	56	3	50	0	49	37	2
17	0	39	36	1	34	0	0	11	2					

尺度水準

- 尺度の4つの種類
 - 比率尺度
 - 間隔尺度
 - 順序尺度
 - 名義尺度

データ行列



有名人データ

<http://note.chiebukuro.yahoo.co.jp/detail/n5672>

	name	height	weight	gender	genre
1	香取慎吾	180.0	71	M	歌手
2	稲垣吾郎	176.0	63	M	歌手
3	木村拓哉	176.0	74	M	歌手
4	草彅剛	170.0	60	M	歌手
5	中居正広	165.0	55	M	歌手
10	大野智	166.0	52	M	歌手
92	木村カエラ	154.0	45	F	歌手
93	米良美一	142.7	35	F	歌手
94	中村昌也	192.0	80	M	俳優
95	阿部寛	189.0	75	M	俳優
96	城田優	188.0	96	M	俳優
97	松重豊	188.0	72	M	俳優
367	ザ・たっち	151.0	43	M	その他
368	池乃めだか	150.0	35	M	その他
369	猫ひろし	147.0	35	M	その他
370	山崎静代(南海)	182.0	64	F	その他
371	富永愛	179.0	70	F	その他
372	熊井友理奈	176.0	59	F	その他

データの記述、要約

●データのグラフ（視覚化）

- 幹葉表示 stem and leaf
- 度数分布図、ヒストグラム histogram
- 箱ひげ図 box and whisker plot

●データの数値的要約

- データの中心の位置
- データの散らばりの大きさ
- データの形、位置

離散データの度数分布

```
>table(gender)
gender
  F   M
203 261
```

```
>table(genre)
genre
その他 歌手 俳優
  234   93  137
```


連続データの度数分布

21

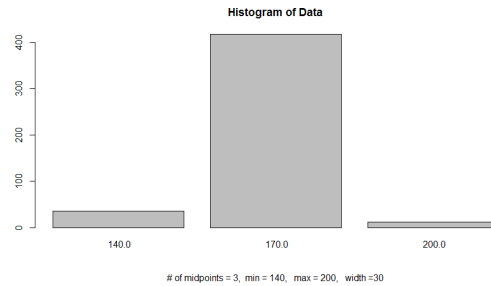
Frequency Distribution

n = 464, nobs = 464

npoints = 3, min = 140, max = 200, width = 30

Categorization of Variable: frequency

NA	<	140.0	<=	155.0	35
155.0	<	170.0	<=	185.0	417
185.0	<	200.0	<=	NA	12



連続データの度数分布

22

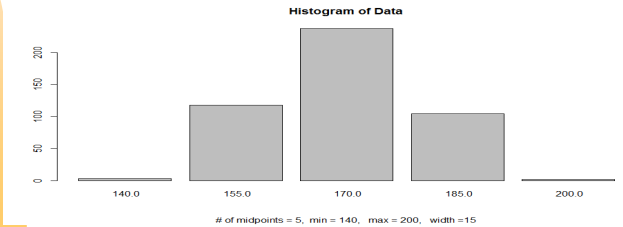
Frequency Distribution

n = 464, nobs = 464

npoints = 5, min = 140, max = 200, width = 15

Categorization of Variable: frequency

NA	<	140.0	<=	147.5	3
147.5	<	155.0	<=	162.5	118
162.5	<	170.0	<=	177.5	237
177.5	<	185.0	<=	192.5	104
192.5	<	200.0	<=	NA	2



連続データの度数分布

23

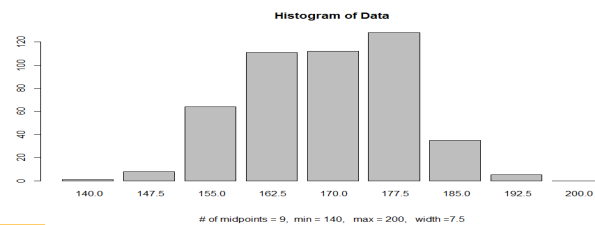
Frequency Distribution

n = 464, nobs = 464

npoints = 9, min = 140, max = 200, width = 7.5

Categorization of Variable: frequency

NA	<	140.0	<=	143.8	1
143.8	<	147.5	<=	151.2	8
151.2	<	155.0	<=	158.8	64
158.8	<	162.5	<=	166.2	111
166.2	<	170.0	<=	173.8	112
173.8	<	177.5	<=	181.2	128
181.2	<	185.0	<=	188.8	35
188.8	<	192.5	<=	196.2	5
196.2	<	200.0	<=	NA	0



ヒストグラム：度数分布

24

147.5 < 155.0 <= 162.5 118

break	midpoint	break	frequency
境界値	代表値	境界値	度数

$$b_{j-1} < x_j \leq b_j \quad f_j$$

度数分布、相対度数分布

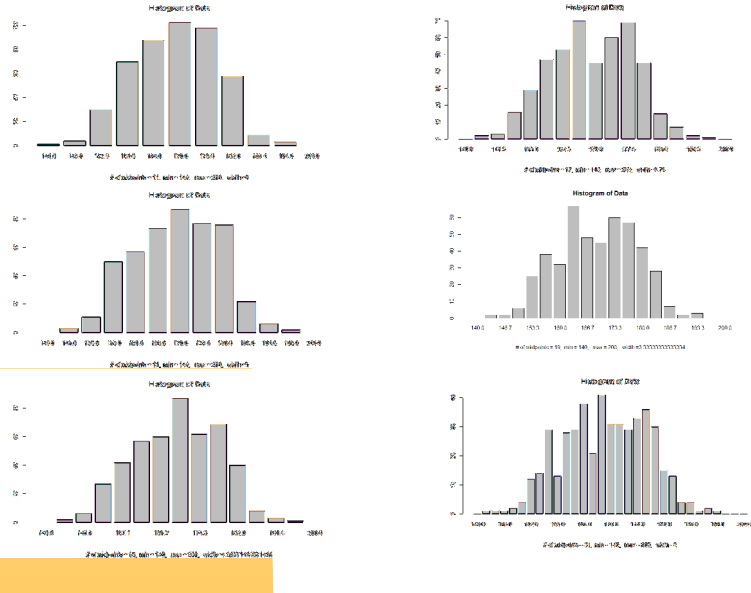
$$(x_j, f_j), \quad j = 1, 2, \dots, k$$

$$(x_j, p_j), \quad j = 1, 2, \dots, k$$

$$\sum_{j=1}^k f_j = n$$

$$p_j = \frac{1}{n} f_j, \quad j = 1, 2, \dots, k \quad \sum_{j=1}^k p_j = 1$$

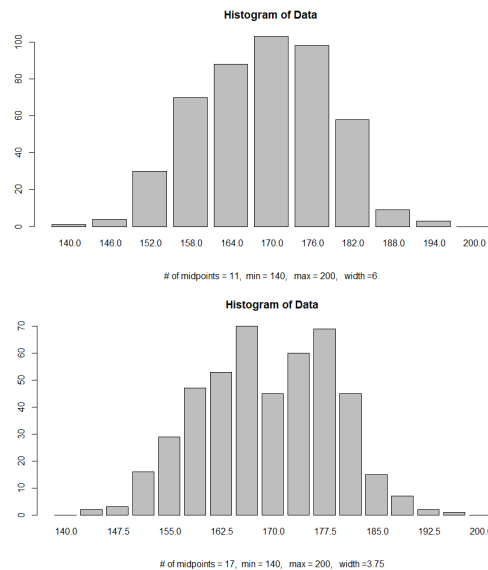
ヒストグラム



ヒストグラム

- 相対度数を用いたヒストグラムは母集団における連続的確率変数の分布を推定している。
- 標本数 n を多く取り、カテゴリ数を増やしていけば、ヒストグラムは母集団分布に近づくはず。
- 最適なカテゴリ数は？
- ちょっと冗長！？

データの中心の位置



一山(単峰)の分布

二山(双峰)の分布

データの中心の位置

- 最頻値 mode
 - 頻度が最も多いカテゴリの値
- 中央値 median
 - 全データを並べ替えた時の真ん中の値
- 平均値 mean
 - ヒストグラムを下から支えてバランスが取れる点

点の選択に伴う損失

29

● 中央値

$$loss = \sum_{i=1}^n (y_i - c)^2 \approx \sum_{j=1}^k f_j (x_j - c)^2$$

● 平均値

$$loss = \sum_{i=1}^n |y_i - c| \approx \sum_{j=1}^k f_j |x_j - c|$$

平均値の計算

30

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \approx \frac{1}{n} \sum_{j=1}^k f_j x_j = \sum_{j=1}^k p_j x_j$$

$$\sum_{j=1}^k f_j = n$$

$$\sum_{j=1}^k p_j = 1$$

分布関数

31

● 相対累積分布関数

■ x_k より小さいデータの割合

$$F_k = \sum_{j=1}^k f_j \quad P_k = \sum_{j=1}^k p_j$$

$$0 \leq P_k \leq 1$$

カテゴリの代表値の関数としての分布関数

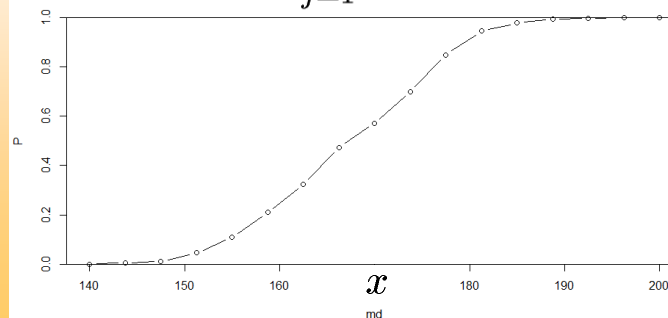
$$P(x_k) = \sum_{j=1}^k p_j$$

分布関数

32

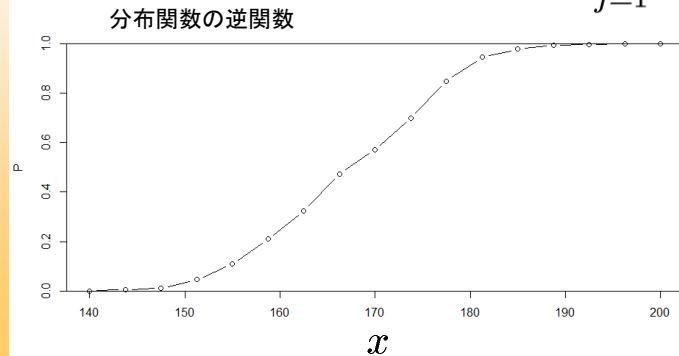
データの値（カテゴリの代表値）を与えると比率を返す関数

$$P(x_k) = \sum_{j=1}^k p_j$$



比率を与えるとデータの値を返す関数

$$x = P^{-1}(p) \qquad P(x_k) = \sum_{j=1}^k p_j$$

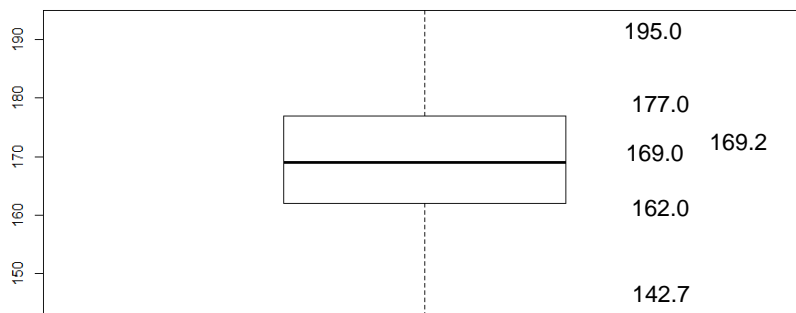


- データの分布の中で、下から $100 \times p \%$ の位置にあるデータの値

- p の値としては 0.10, 0.25, 0.50, 0.75, 0.90

- 四分位値 quartile

$$Q_1 = P^{-1}(0.25), \quad Q_2 = P^{-1}(0.50), \quad Q_3 = P^{-1}(0.75)$$



標準偏差 = 9.48, IQR = 15, 四分位偏差 = IQR/2 = 7.5

- 範囲 レンジ range
 - 最大値 - 最小値

- 四分位範囲 Inter Quartile Range IQR
 - Q3-Q1

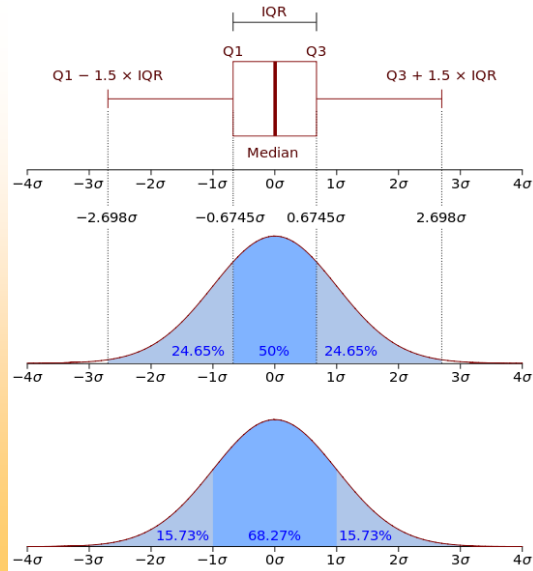
- 標準偏差 = 分散 S^2 の平方根 S

$$S^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \approx \sum_{j=1}^k p_k (x_k - \bar{y})^2$$

$var(y)$

variance, standard deviation, sd or std

大体の目安



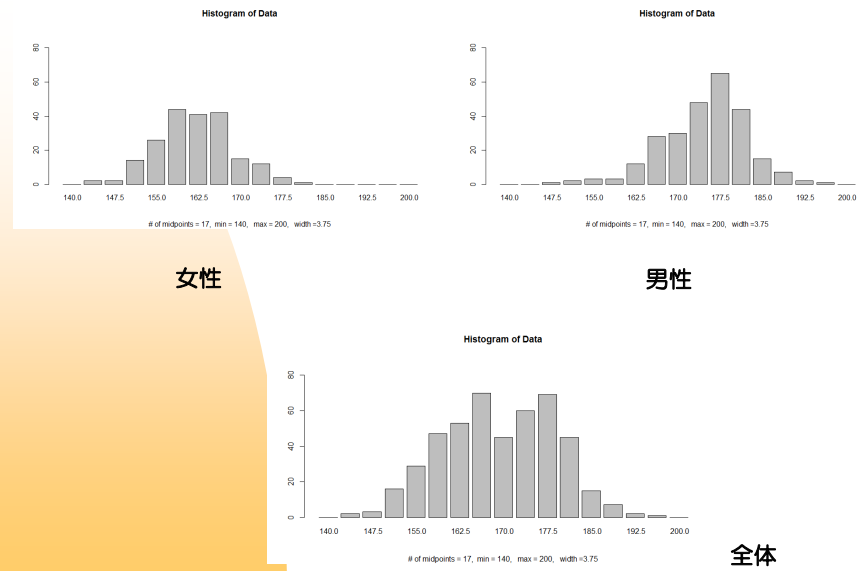
中央値 ± 0.5 IQR = 50%

平均 ± 1 標準偏差 $\approx 68\%$

平均 ± 2 標準偏差 $\approx 95\%$

平均 ± 3 標準偏差 $\approx 99.7\%$

層別の分布

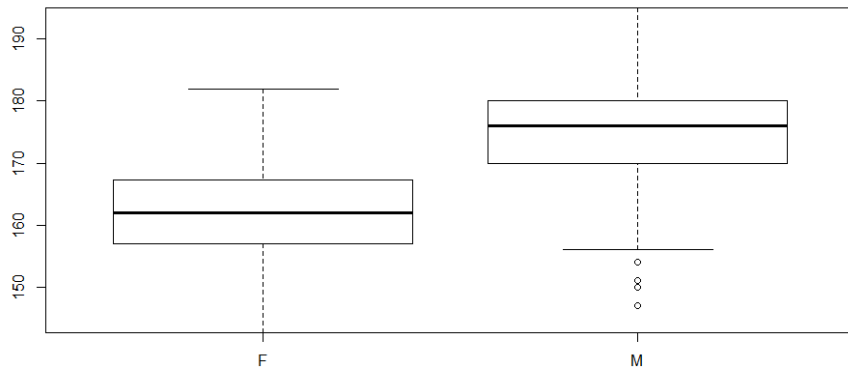


女性

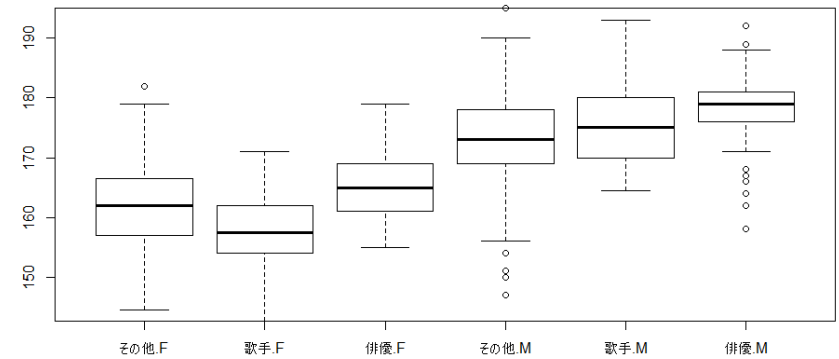
男性

全体

層別の分布

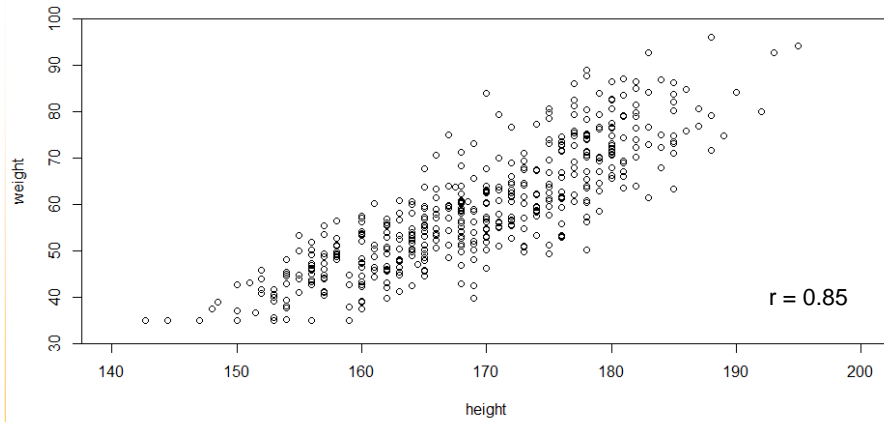


層別の分布



二変数のプロット：散布図

41



相関係数

42

● 共分散

$$C_{y,z} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})$$

● 相関係数

共分散をそれぞれの標準偏差で除したもの

$$r_{y,z} = \frac{C_{y,z}}{\sqrt{\text{var}(y)\text{var}(z)}} \quad 0 \leq r_{y,z} \leq 1$$

層別の散布図

43

