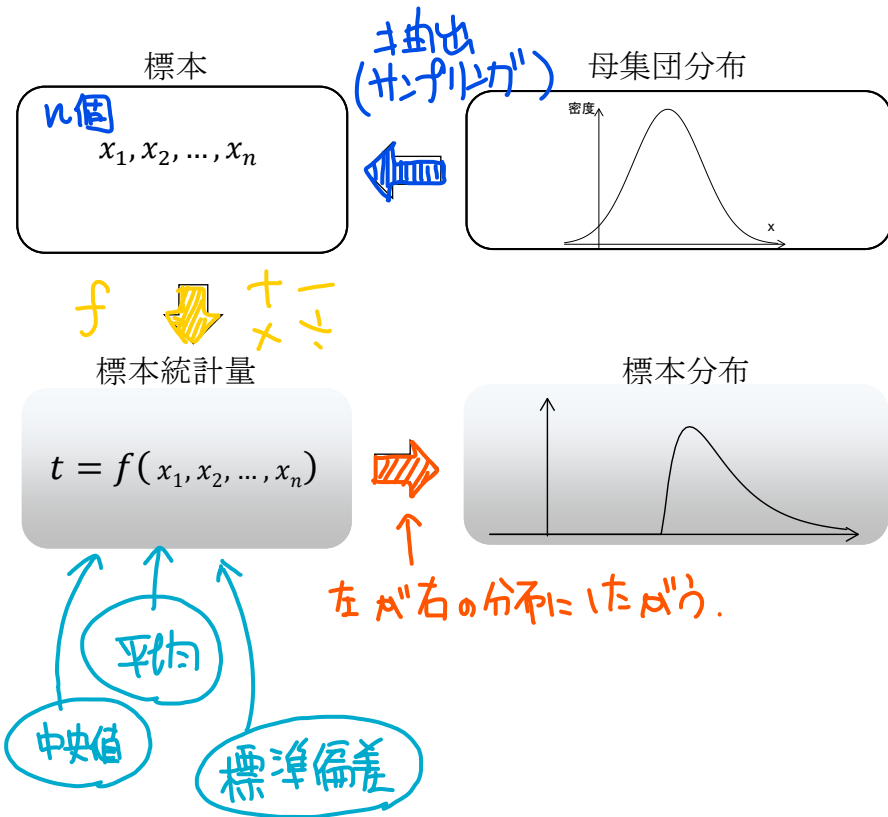


# 見取り図





(1) 標本 Sample

得られたデータのことを標本と呼ぶ。多くの場合全てを報告するのではなく何らかの指標で標本全体を代表させる。

(2) 母集団 Population

研究者が想定する「考えられる対象をすべて含む集団」。  
 ※標本が抽出される際に、確率的な揺らぎが生じる。

サンプル  
 標本(観測データ)を + - x ÷ して  
 つくった量のこと。

④ 関数 f

これは、入力値が  
 どの様に出力値に  
 変換されるかを示した  
 もののこと。

$$f(x) = 2 \times x$$

↑ 入力値      ↑ 出力値

$$f(x) = 5 + x$$

↑ データという入力値      ↓ 統計量

👉 (標本) 統計量

統計量とは、ある任意の関数  $f$  が標本  $x_1, x_2, \dots, x_n$  を引数にとったときに返す返り値のことである。

$$T = f(x_1, x_2, \dots, x_n)$$

↑ データ      ↓ 入力値

↑ 統計量      ↓ 変換      ↓ 入力値

標本から計算される指標を(標本)統計量と呼び、とりわけデータをうまく要約している統計量を代表値と呼ぶ。

例：平均

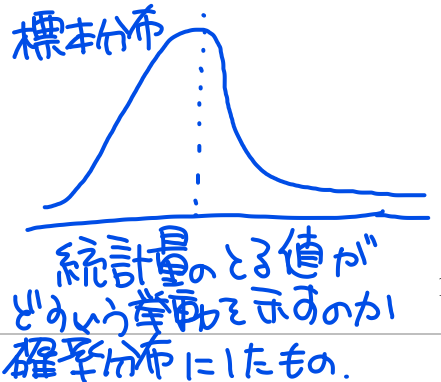
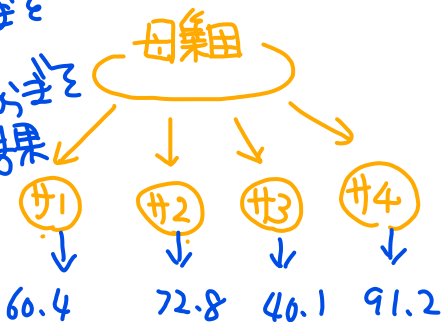
$$\text{平均} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

これは、データに + と ÷ をあてがって新しい量を生成したもの。

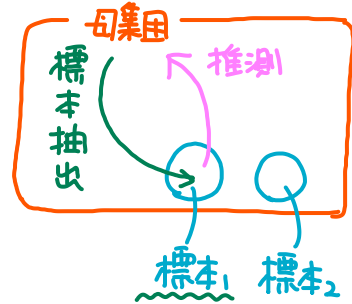
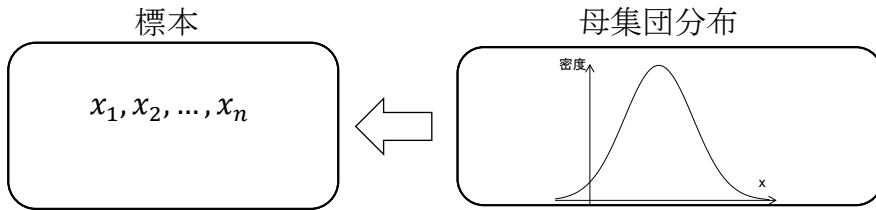
(4) 標本分布 Sampling distribution

母集団から標本がサンプリングされる際の確率的な揺らぎのせいで、統計量は一つの値には決まらない。その結果統計量が見せる確率的挙動を表したものの。

← サンプリングというプロセスが確率的ゆらぎを招くが、結果的に統計量も確率的ゆらぎを内包にしている。分岐が主たる意味



④ 母集団と標本



(1) 標本とヒストグラム

① 標本 (サンプル)

これは、研究者が観測するデータのこと。

② サンプルサイズ

これは、標本の中に含まれている観測値の数のこと。

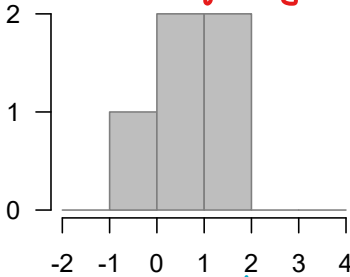
サンプルサイズは3.

ID	容
1	5
2	5
3	5

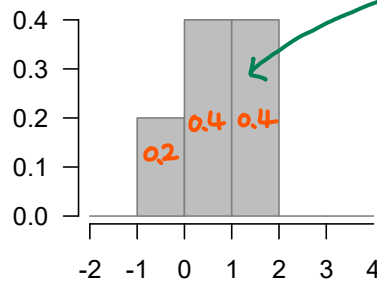
③ ヒストグラム (柱状図)

これは、設定された各区間にどのくらいの頻度で観測値が存在しているのかを柱で表したグラフ。

(1) 粗頻度: 回数をy軸  
Raw Frequency



(2) 相対頻度: 比率をy軸



④ 面積が割合!

灰色の面積の合計 = 1 (100%)

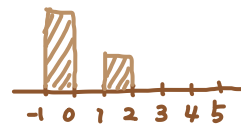
1~2の区間 = 2回 → 2/5 = 0.4  
 0~1 " 2回 → 2/5 = 0.4  
 -1~0 " 1回 → 1/5 = 0.2  
 ⇒ 全区間 5回

④ 棒グラフとの違い

(1) ヒストグラム

質問 ヒストグラムと棒グラフは同じですか?

違います。ヒストグラムは「区間」に対してその頻度を長さで表現していますが、棒グラフは「点(値)」に対してその大きさを長さとして表しています。



区間に対して定まる

(2) 棒グラフ



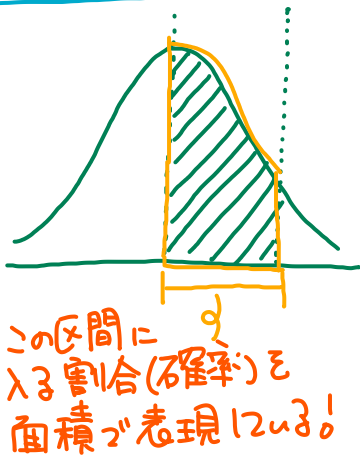
点(値)に対して定まる

(2) 母集団分布

① 母集団

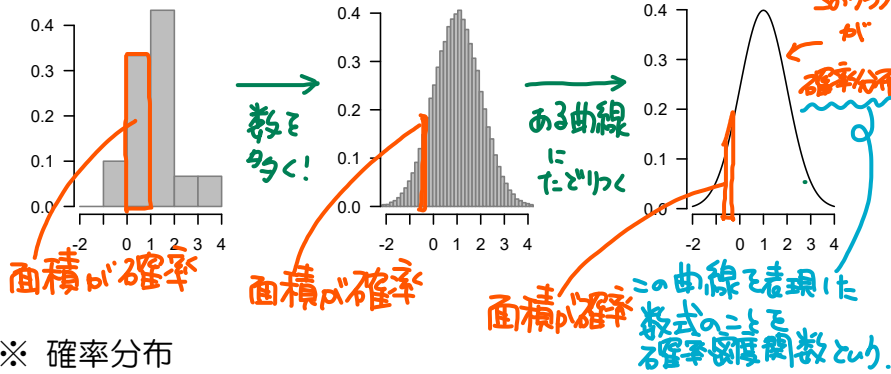
研究者が想定する「考えられる対象をすべて含む集団」。  
 ※標本が抽出される際に、確率的な揺らぎが生じる。

① 確率密度関数の見方



② 母集団分布

標本のヒストグラムの行きつく先。母集団の中心やばらつきが表されている。



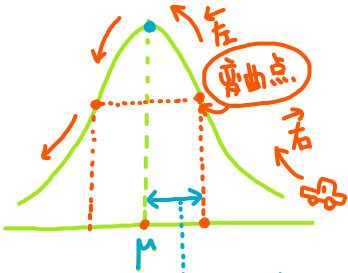
※ 確率分布

「区間」に対応する確率を示しているグラフのこと。

③ (確率) 密度関数 Probability density function

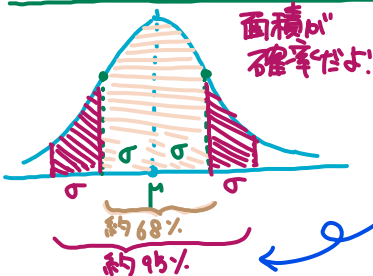
確率分布に示された曲線を表すグラフの式。

① 変曲点と  $\sigma$



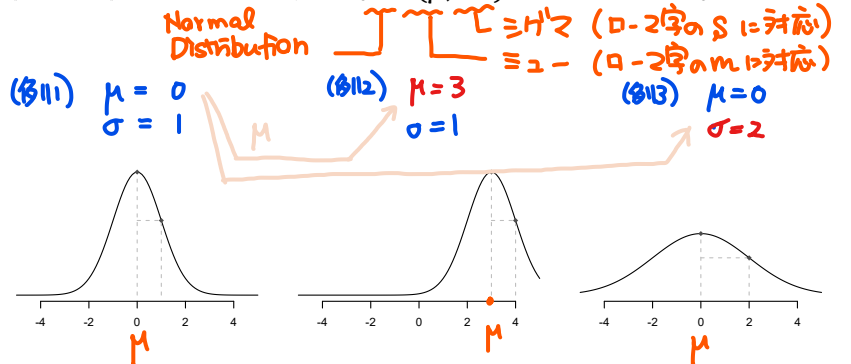
- ①  $\mu$ : 位置パラメータ (分布の中心)
- ②  $\sigma$ : 尺度 (分布の広がり)

① 正規分布と確率



(3) 正規分布

ランダムな誤差が積み重なると生まれる左右対称の釣鐘状の形状をした確率分布。 $N(\mu, \sigma^2)$ のように表す。



- (特徴1) 数学的に扱いやすい!
- (特徴2) 無標な分布

① より厳密には、 $1.96 \times \sigma$  の区間が全体の95%に対応してる。

(1) 統計量と標本分布

① **統計量** Statistic

これは、標本に基づいて計算された量。

② **標本分布** Sampling Distribution

これは、統計量が描くヒストグラムの行き着く先。

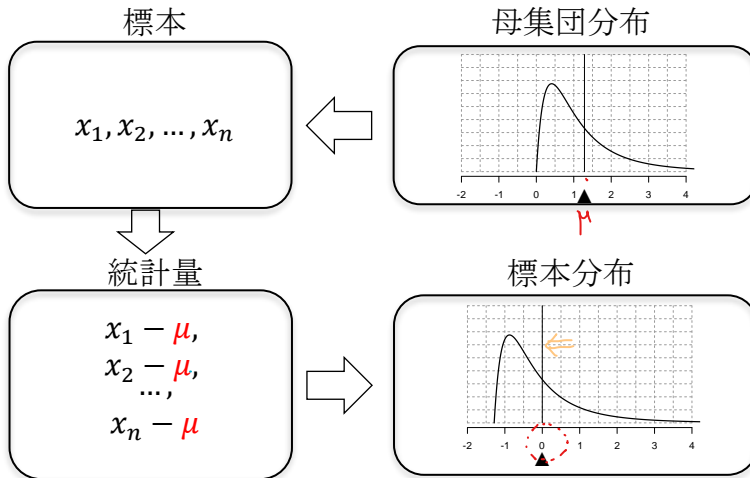
(2) 統計量と標本分布の具体例

① **操作1**：中心化

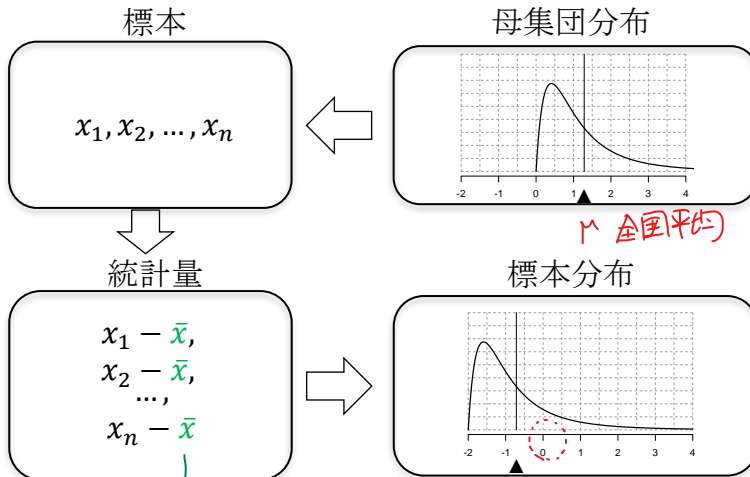
平均を引き、データの中心を動かす平行移動のこと。

データ - 中心

(A) **全知全能の視点**



(B) **人間の視点**



④ **具体例**

$x_1$  : 出席番号1の人の点数

$x_2$  : " 2 "

⋮

$x_n$  : " n "

標本のサンプルサイズ

④ **記法** :  $\bar{\quad}$  (バー)

ある変数の平均を  
 $\bar{\quad}$  (バー) をつけて表す

(例)  $x_1, x_2, \dots, x_n$   
 $\downarrow$   
 $\bar{x}$  : エックス・バー

① 割り算をすること

これは、分母にきている「単位」何個分なのかを計算すること。

120個のアイスクリーム

1ケース12個はいり

= 10ケース分の量がある

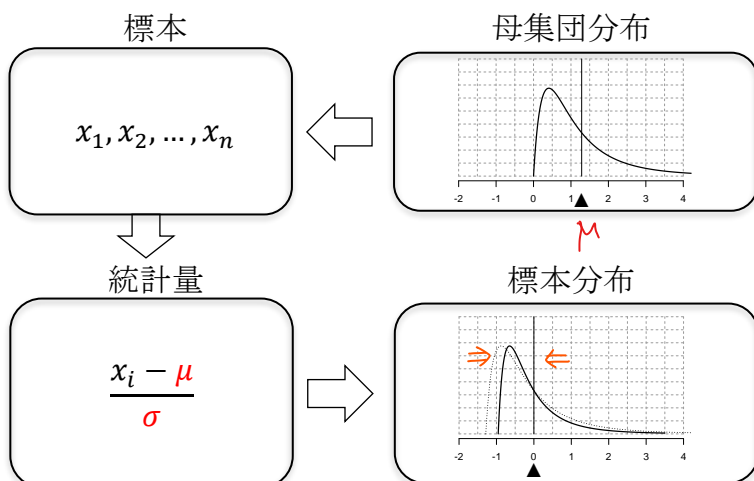
② **操作2**：標準化

平均を引き、データの中心を動かす平行移動をした値を、ばらつきの尺度である標準偏差で割ること。

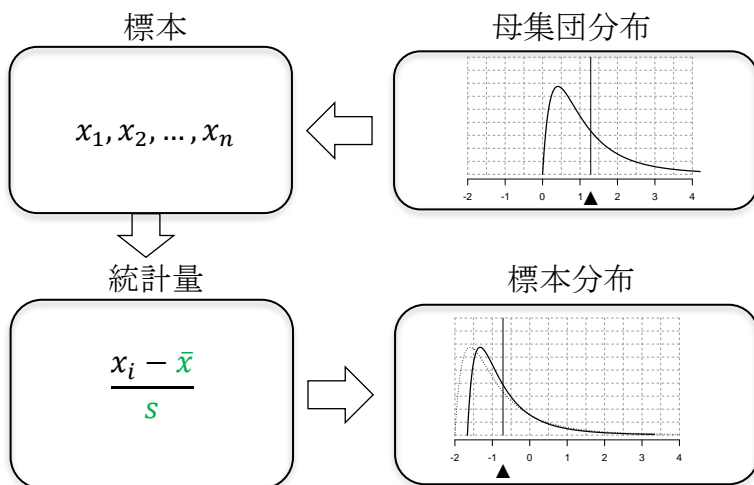
データ - 中心

ばらつきの尺度

(A) **全知全能の視点**



(B) **人間の視点**

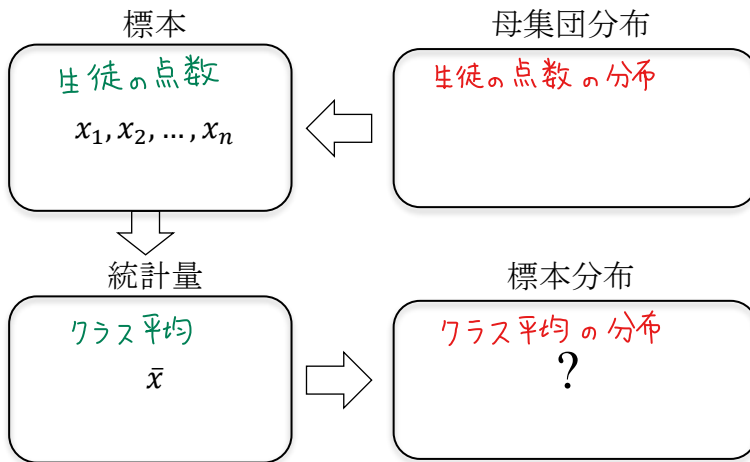


(3) **中心極限定理** Central Limit Theorem

これは、平均という統計量がしたがう標本分布が正規分布になるという定理。母集団の形状で二種類のものがある。

① 具体例

「大学入試共通テスト」を受けた*i*番目の生徒の英語の点を $x_i$ とする。このときクラス平均がしたがう分布は？

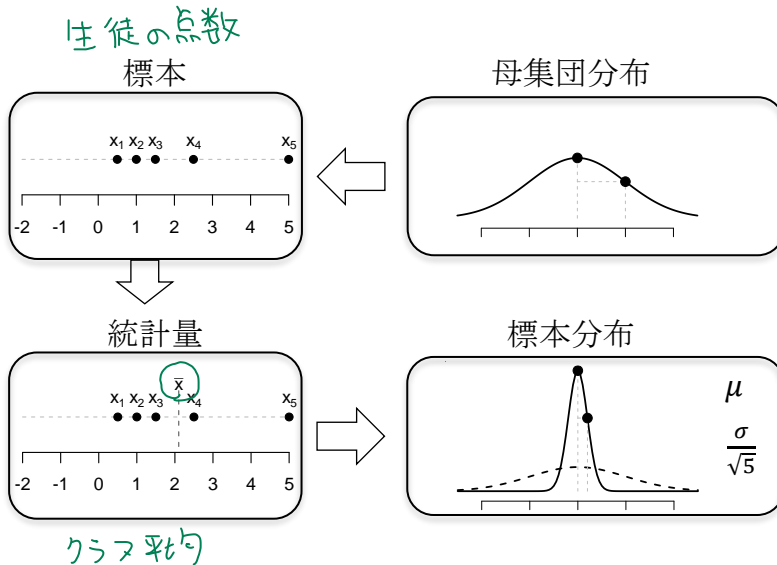


② **バージョン1**：母集団分布が**正規分布**のとき

任意の母集団では、サンプルサイズ  $n$  が大きくなっても、 $\bar{x}$ の標本分布は $N(\mu, \frac{\sigma^2}{n})$ になる。

**サンプルサイズ  $n$  が小さいケース**

$n = 5$



④ サンプルサイズ  $n=1$

これは、すべてのクラスサイズが「マンツーマン」という個別指導の状況を表わしている。

$$\bar{x} = x_1$$

$$x_1 \sim N(\mu, \sigma^2)$$

$$\bar{x} \sim N(\mu, \sigma^2)$$



## 再生性 (正規分布)

(ケース 1) 別々の正規分布に従う場合

$$x_{1j} \sim N(\mu_1, \sigma_1^2) \quad \text{英語の点}$$

$$x_{2j} \sim N(\mu_2, \sigma_2^2) \quad \text{数学の点}$$

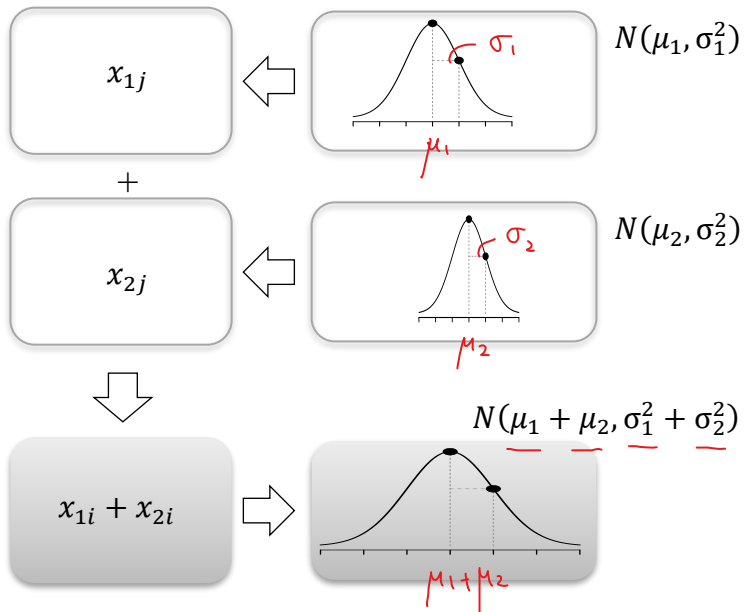
⋮

$$x_{nj} \sim N(\mu_n, \sigma_n^2) \quad \text{美術の点}$$

+) \_\_\_\_\_

$$x_{1j} + x_{2j} + \dots + x_{nj} \sim N(\mu_1 + \mu_2 + \dots + \mu_n, \sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2)$$

総合点



(ケース 2) 同一の正規分布に従う場合

$$x_{1j} \sim N(\mu, \sigma^2)$$

$$x_{2j} \sim N(\mu, \sigma^2)$$

⋮

$$x_{nj} \sim N(\mu, \sigma^2)$$

+) \_\_\_\_\_

$$x_{1j} + x_{2j} + \dots + x_{nj} \sim N(n\mu, n\sigma^2)$$

↓ サンプルサイズ  $n$  で割る

$$\frac{1}{n}(x_{1j} + x_{2j} + \dots + x_{nj}) \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

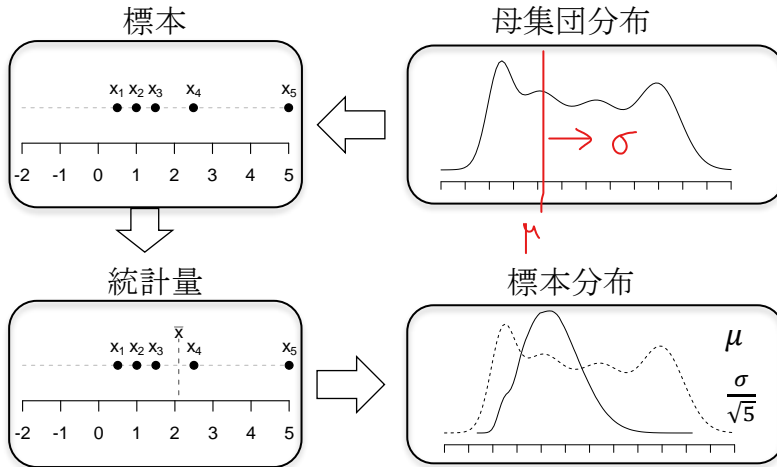


③ バージョン2：母集団分布が正規分布以外の場合

任意の母集団では、サンプルサイズ  $n$  が大きければ大きいほど、 $\bar{x}$  の標本分布は  $N\left(\mu, \frac{\sigma^2}{n}\right)$  に近づいていく。

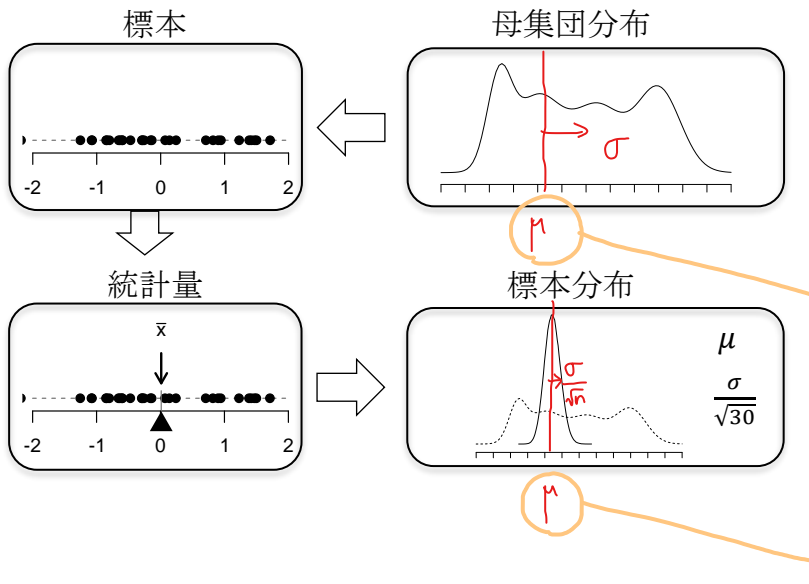
サンプルサイズ  $n$  が小さいとき

$n = 5$



サンプルサイズ  $n$  が大きくなってきたとき

$n = 30$





## 数学が苦手な人のために

( 英語 ) 単語・句・節

S V  
I remember

○  
↑  
that

単語 (名詞)

visiting my friends 句 (名詞句)

that I visited my friends 節 (名詞節)

( 数学 ) 数式の中における埋め込み

ここに入る表現が複雑化する

$$\frac{\boxed{\text{データ}} - \boxed{\text{中心}}}{\boxed{\text{ばらつきの尺度}}}$$



## この後何度も使うことになる操作

(1) 標準化

データから中心を引いて、それをばらつきの尺度で割る

(2) 中心極限定理

母集団が正規分布だと、平均の標本分布は正規分布になる

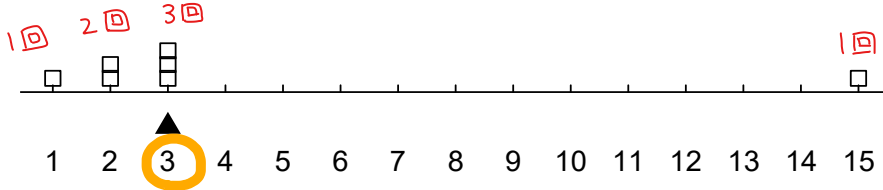
(3) (正規分布の) 再生性

正規分布に従う変数を足した和も (差も) 正規分布になる

■ 基本的な統計量I (データの中心)

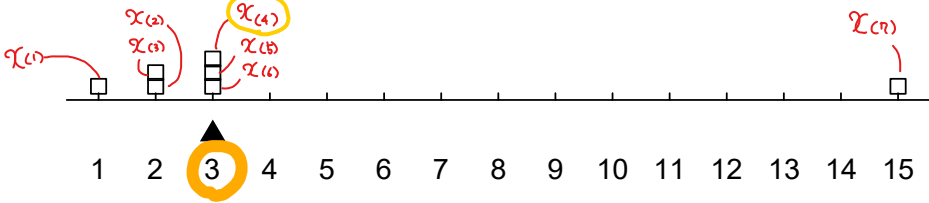
(1) (標本) **最頻値** Sample Mode

これは、**頻度** という観点からデータの真ん中を定めたもので、得られたデータの中で最も多く得られた値のこと。



(2) (標本) **中央値** Sample Median

これは、**順位** という観点からデータの真ん中を定めたもので、小さい順に並べた時の真ん中にくる値のこと。



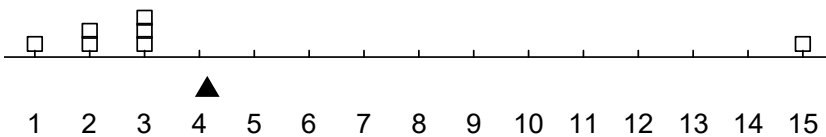
$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(\frac{n}{2})} \leq x_{(\frac{n}{2}+1)} \leq \dots \leq x_{(n)}$$

最小値  最大値

$$MD = \begin{cases} x_{(\frac{n+1}{2})} & n \text{ が 奇数のとき} \\ \frac{1}{2} (x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}) & n \text{ が 偶数のとき} \end{cases}$$

(3) (標本) **平均** Sample Mean

これは、**バランス** という観点から真ん中を定めたもので、値をすべて足して総数で割った値のこと。



$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$= \frac{1}{n} (x_1 + x_2 + \dots + x_i + \dots + x_n)$$

① **代表値**

これは、データ (= 標本) を 1 つの値に代表させて示すときに用いる統計量のこと。

- ① データの中心
- ② データのばらつき
- ③ データ間の関連度

② **順序統計量**

これは、順位による定義される統計量のこと。

- (例)
- 最大値
  - 中央値
  - 最小値

③ **記法** :  $x_1$  と  $x_{(i)}$

① 生のデータ

出席番号	英語
$x_1 \rightarrow 1$	60
$x_2 \rightarrow 2$	80
$\vdots$	$\vdots$
$\vdots$	26
$\vdots$	$\vdots$

② 並べかえたデータ

出席番号	英語
$x_{(3)} \rightarrow 3$	26
$x_{(6)} \rightarrow 17$	35
$\vdots$	36
$\vdots$	$\vdots$

④ **総和記号**  $\sum_{i=1}^n x_i$

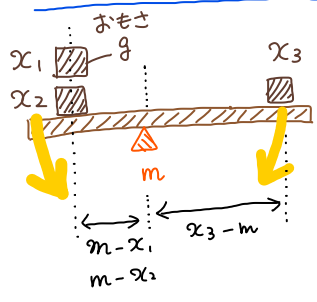
これは Σ (シグマ) と読み、英語の S に対応するもの。  
(Σ: 大文字, σ: 小文字)

the sum of  $x_i$  from  $i=1$  to  $n$ .

cf. 総積記号  $\prod_{i=1}^n x_i$

$$\prod_{i=1}^n x_i = x_1 \times x_2 \times \dots \times x_n$$

① 平均が  $\frac{1}{n} \sum_{i=1}^n x_i$  の理由



$$g \times (m-x_1) = g \times (x_3-m)$$

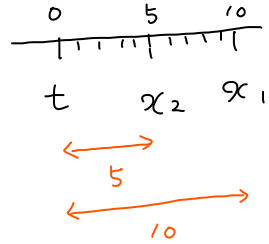
$$+ g \times (m-x_2)$$

$$3m = x_1 + x_2 + x_3$$

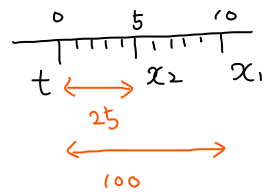
$$m = \frac{x_1 + x_2 + x_3}{3}$$

② "距離" の考え方を拡張

① ユークリッド距離



② 平方ユークリッド距離

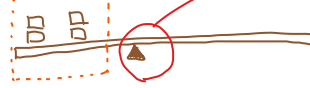


どちらのアプローチをとっても  $x_1$  が  $x_2$  より  $t$  から遠いという点は捉えらる。  
 → と2つ②のケースも、その「遠さ」を測れる指標として採用してもいいじゃん。

③ 外れ値 (outlier)



データが集まっているところ  
 → と3  
 平均がデータが分布している中心部にはない



外れ値: データが分布しているところから離れたところに位置する点のこと。

データとの距離の最小化

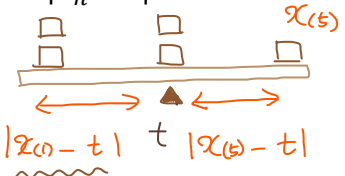
① 中央値: 一次の距離の最小化

「 $t$  と各点  $x_i$  の 差の絶対値 の総和を最小化する」という基準を満たすデータの中心。  
 (正規が測った距離)   
 (ユークリッド距離)   
 (=我々が日常的に使っている距離)

⇒ すなわち、次の基準  $T_1$  を最小化する  $t$  が中央値。

$$T_1 = |x_1 - t| + |x_2 - t| + \dots + |x_n - t|$$

$$= \sum_{i=1}^n |x_i - t|$$



↑ 絶対値を計算するとき場合分けが必要 (=手間)

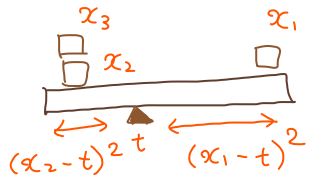
② 平均: 二次の距離の最小化

「 $t$  と各点  $x_i$  の 差の二乗 の総和を最小化する」という 最小二乗基準 を満たすデータの中心。

⇒ すなわち、次の基準  $T_2$  を最小化する  $t$  が平均。

$$T_2 = (x_1 - t)^2 + (x_2 - t)^2 + \dots + (x_n - t)^2$$

$$= \sum_{i=1}^n (x_i - t)^2$$



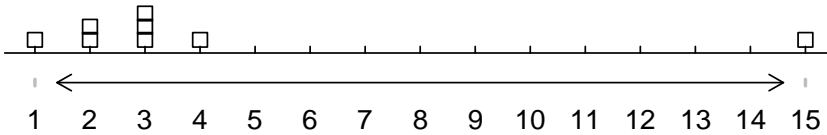
それぞれの指標のメリットとデメリット

- (特徴1) 中央値、平均と違い 最頻値 は存在しないときがある。
- (特徴2) 中央値と違い平均は 外れ値 の影響を受けやすい。
- (特徴3) 中央値と違い平均は 数学的な 取り扱いが楽であり、かつ、いろいろときれいな性質がある。

## ■ 基本的な統計量Ⅱ (1次データのばらつき)

### (1) 標本範囲 (最小統計量と最大統計量) Sample Range

これは、最大値 ( $x_{(n)}$ ) と最小値 ( $x_{(1)}$ ) に注目し、その差を測って得られる値のこと。

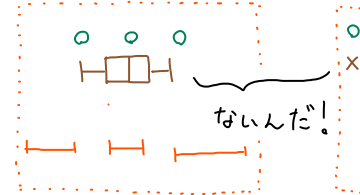


$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

$$Rg = x_{(n)} - x_{(1)}$$

### ④ 箱ひげ図の例外

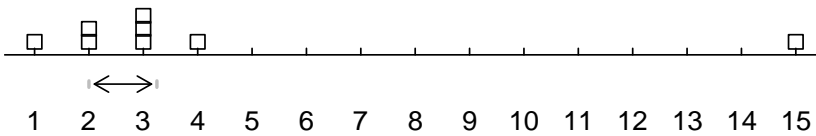
「外れ値 (outlier)」



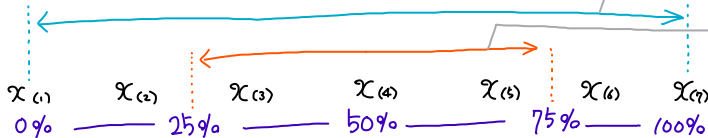
外れ値

### (2) 四分位範囲 Interquartile Range

これは、データを小さい順に並べたときの下位 25% (第 1 四分位点  $Q_1$ )、75% (第 3 四分位点  $Q_3$ ) に位置するデータの距離 (L1 距離) を出したものです。



$$IQR = Q_3 - Q_1$$



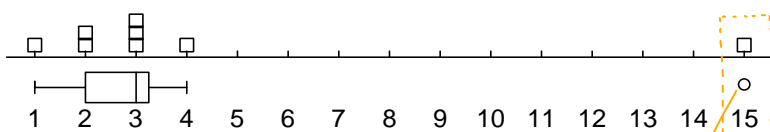
標本範囲 (Range)

最大値 - 最小値

四分位範囲

### ※ 箱ひげ図

これは、最大値、最小値、第 1~3 位四分位点を表す図。外れ値 を示すために、IQR の 1.5 倍の範囲の外にあるものは、独立して表す。



標本範囲

四分位範囲



$Q_0$   $Q_1$   $Q_2$   $Q_3$   $Q_4$   
0% 25% 50% 75% 100%

(中央値)

外れ値

# ① 偏差

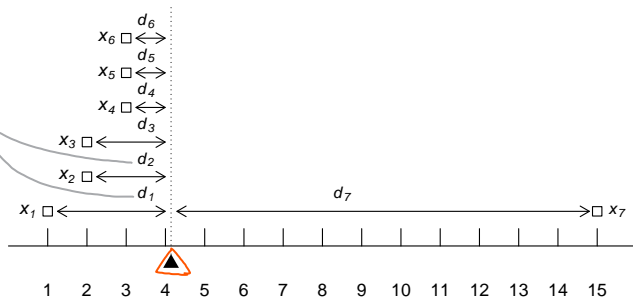
(3) 平均偏差 (一次の距離に基づく指標)

$$|x_1 - \bar{x}| = d_1$$

$$|x_2 - \bar{x}| = d_2$$

⋮

$$|x_7 - \bar{x}| = d_7$$



# ② 偏差を足し合わせる

$$|x_1 - \bar{x}|$$

$$+ |x_2 - \bar{x}|$$

+

$$+ |x_7 - \bar{x}|$$

$$= |x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_7 - \bar{x}|$$

$$= \sum_{i=1}^7 |x_i - \bar{x}|$$

## ① 偏差 Deviation

これは、一次距離 で測った平均からの距離。

$$d_i = x_i - \bar{x}$$

## ② 平均偏差 Mean Deviation

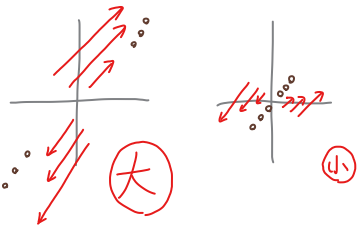
これは、偏差の平均。一次の距離 に基づいた指標。

平均している  $MD = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$  紫色の矢印の長さを全部足しあわせたのがここ

# ③ 分散

(4) 標本分散 (二次の距離に基づく指標)

これは、自分自身との共分散であり、ばらつき  
の指標である



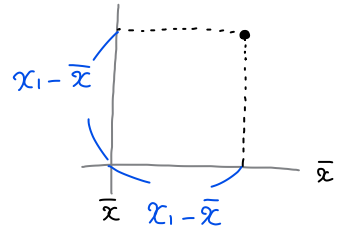
## ① 偏差平方和

x軸 × y軸

これは、二次距離 で測った平均からの距離。

$$d_i^2 = (x_i - \bar{x})^2$$

$$SSD = \sum_{i=1}^n (x_i - \bar{x})^2$$



## ② 分散 Variance

これは、偏差平方和の平均。

二乗

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

面積

## ③ 標準偏差 Standard Deviation

これは、分散の平方根。

√をとります

$$s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

長さ

# ④ 変数 x の分散の表記

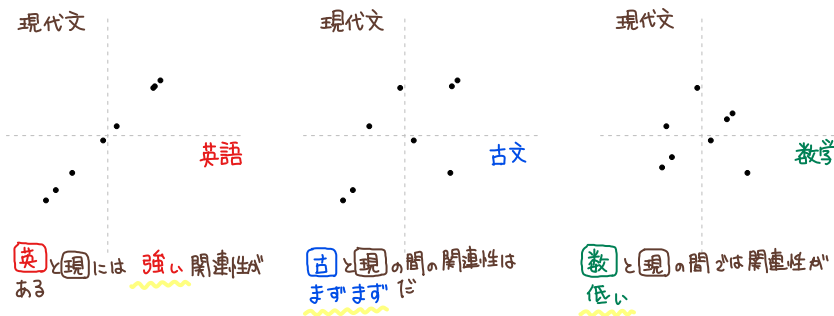
$$S^2 \text{ (文脈上明白)}$$

あるいは

$$S_x^2 \text{ (文脈上変数が2つ以上あって、どけを指すのかはきりさせたいとき)}$$

## ■ 基本的な統計量 III (二つの変数の 相関)

【目的】x軸とy軸の関連性の方向と大きさを測る指標を作りたい。



### (1) 共分散

これは、平均と各点との間の二次の距離(面積)の平均。

$$s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

※共分散には最大値と最小値が存在する。

$$-s_x s_y \leq s_{xy} \leq s_x s_y$$

### (2) 相関係数 (ピアソンの積率相関係数)

これは、共分散が最大値と比較したときに占める割合。

$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

今回とった共分散

関連度がMAXのときの共分散

※相関係数には最大値と最小値が存在する。

$$-1 \leq r_{xy} \leq 1$$

↑ -100%

↑ 100%

### ① 最大値が割ること

Aさん: 90点 (990点満点)  
Bさん: 90点 (200点満点)

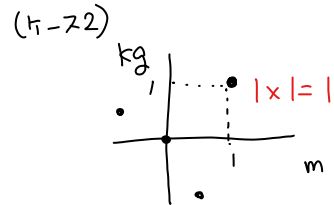
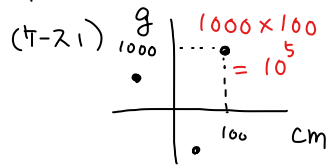
$$\frac{90}{990} \doteq 0.09 (9\%)$$

$$\frac{90}{200} = 0.45 (45\%)$$

最大値を分母

### ① 共分散の弱点

単位に依存してしまう。



### ① 共分散の上限と下限

共分散  $s_{xy}$  は、関連性が最大化するときに、次の値をとる

$$s_{xy} = s_x s_y$$

(右肩あがり)

$$s_{xy} = -s_x s_y$$

(右肩下がり)

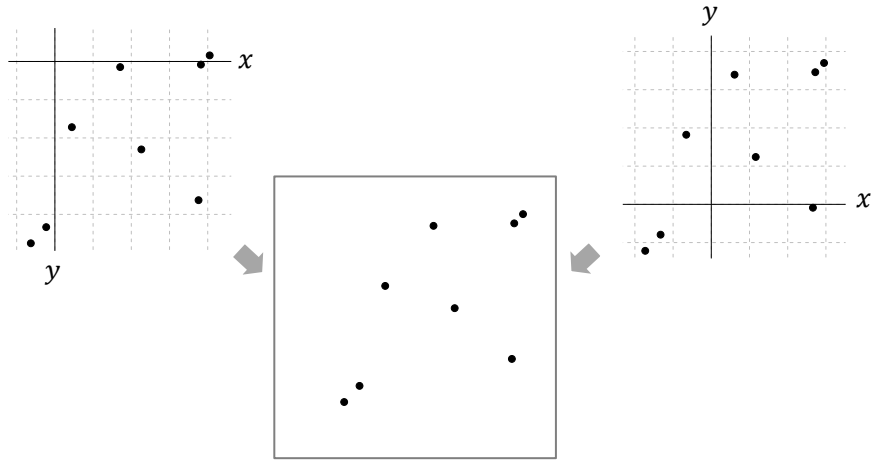
### ① 相関係数の良いところ

- ① 正負: 関連の方向
- ② 値: 関連度の強さ  
1: 相関が(正の方向に)MAX  
0: 相関がない  
-1: 相関が(負の方向に)MAX
- ③ 場合分け不用
- ④ サンプルサイズに依存しない
- ⑤ データの単位に依存しない

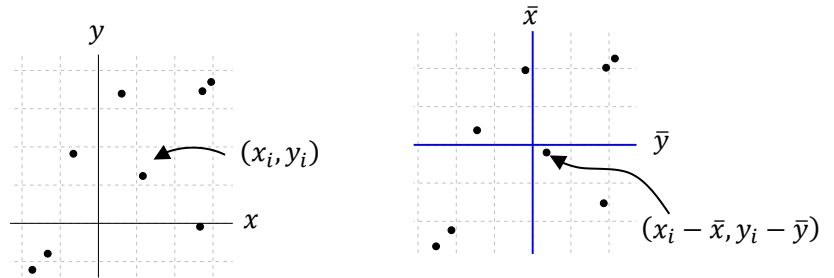


### (準備) 軸を取り払う

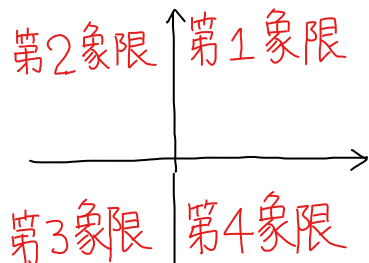
もともとの  $x$  軸、 $y$  軸の位置は関連度合いには無関係なので、取っ払ってしまいたい。



しかし、座標自体は維持したいので、新しく、横軸を  $\bar{y}$ 、縦軸を  $\bar{x}$  の位置に据えた新しい座標を考える。これはもともとの座標を平行移動したことに相当する。



### ① 象限



### ① 第1案の良いとこ3

① 正負：関連の方向

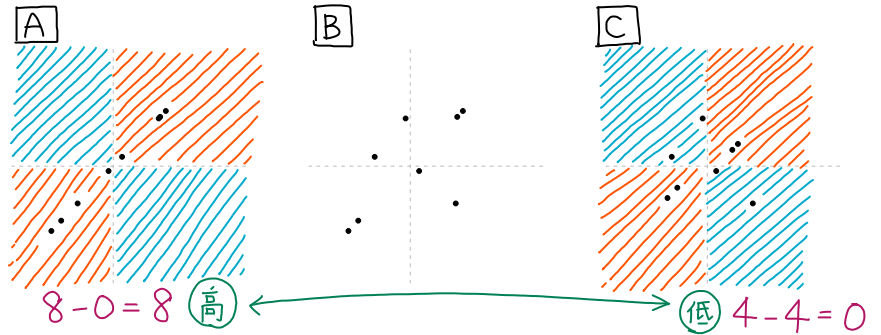
- 正：↗
- 0：↔
- 負：↘

② 値の大きさ：関連の強さ



### (第一案) 球の個数の差を考える

(式1) 「第一／三象限の球の数」 - 「第二／四象限の球の数」



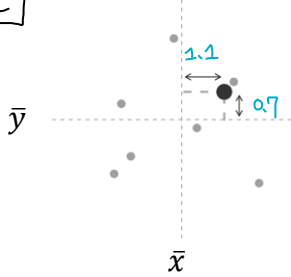
問題点：ケース B とケース C の違いを捉えられない。



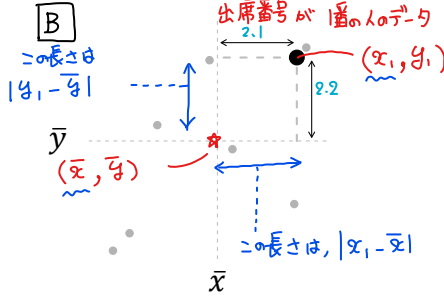


### (第二案) 中心からの離れ具合を「足し算」で捉える

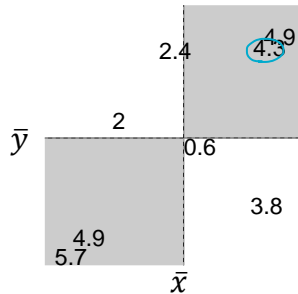
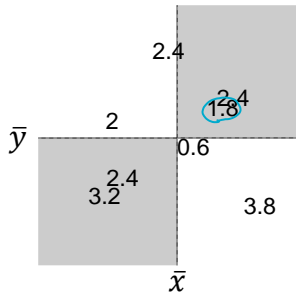
C



B



(式2) 「灰色の領域の数字の和」 - 「白色の領域の数字の和」



問題点：数式に場合分けが存在してしまう。

### ① 第2案の良いところは

- ① 正負：関連性の方向
- ② 値：関連性の強さ
- ③ BとCを区別できる

### ② 数式で表現 (第2案)

① 出席番号1の人の中心からの長さ

$$2.1 + 2.2 = |x_1 - \bar{x}| + |y_1 - \bar{y}|$$

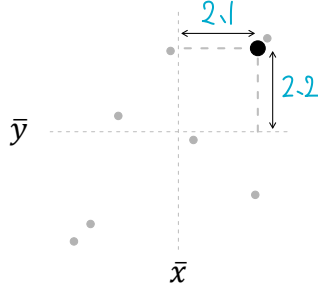
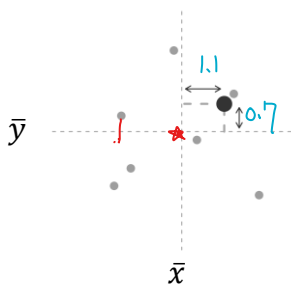
② 全2の人をまとめて、中心からの長さ

$$\begin{aligned} & 4.3 + 4.9 + 4.9 + 5.7 - 2.4 - 2 - 0.6 - 3.8 \\ & + |x_1 - \bar{x}| + |y_1 - \bar{y}| + |x_2 - \bar{x}| + |y_2 - \bar{y}| + |x_3 - \bar{x}| + |y_3 - \bar{y}| + |x_4 - \bar{x}| + |y_4 - \bar{y}| - (|x_5 - \bar{x}| + |y_5 - \bar{y}|) - (|x_6 - \bar{x}| + |y_6 - \bar{y}|) - (|x_7 - \bar{x}| + |y_7 - \bar{y}|) - (|x_8 - \bar{x}| + |y_8 - \bar{y}|) \end{aligned}$$

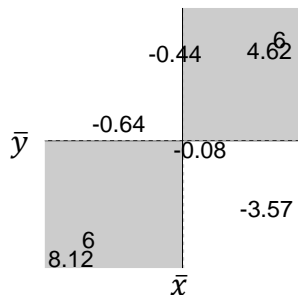
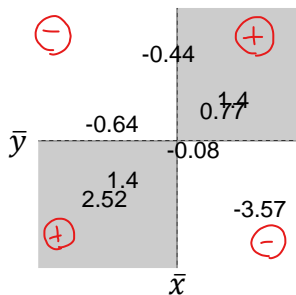
= 11.1



### (第三案) 中心からの離れ具合を「掛け算」で捉える



(式2) 「灰色の領域の数字の和」 - 「白色の領域の数字の和」



問題点：サンプルサイズに依存してしまう。

### ③ 数式で表現 (第3案)

① 出席番号1の人

$$2.1 \times 2.2 = (x_1 - \bar{x})(y_1 - \bar{y})$$

② 全2の人をまとめて

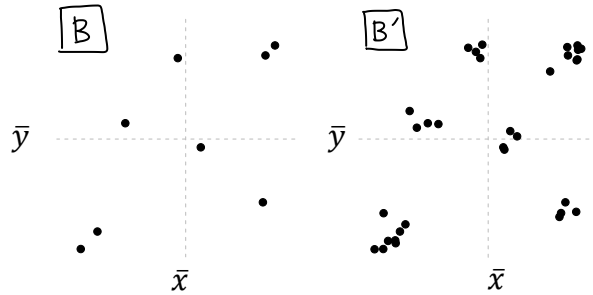
$$\begin{aligned} & 6 + 4.62 + 6 + 8.12 + (-0.44) + (-0.64) + (-0.08) + (-3.57) \\ & + (x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + (x_3 - \bar{x})(y_3 - \bar{y}) + (x_4 - \bar{x})(y_4 - \bar{y}) + (x_5 - \bar{x})(y_5 - \bar{y}) + (x_6 - \bar{x})(y_6 - \bar{y}) + (x_7 - \bar{x})(y_7 - \bar{y}) + (x_8 - \bar{x})(y_8 - \bar{y}) \end{aligned}$$

$$= \sum_{i=1}^8 (x_i - \bar{x})(y_i - \bar{y})$$



### (第四案) サンプルサイズの影響を無くす

採用!



(式3)  $\frac{1}{n} \times (\text{「奇数象限の数字の和」} - \text{「偶数象限の数字の和」})$

共分散

$$s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

準備:  
平行移動を  
する

第4案  
サンプルサイズの影響を取りのぞく

第2案  
中心からの遠さを  
考慮する

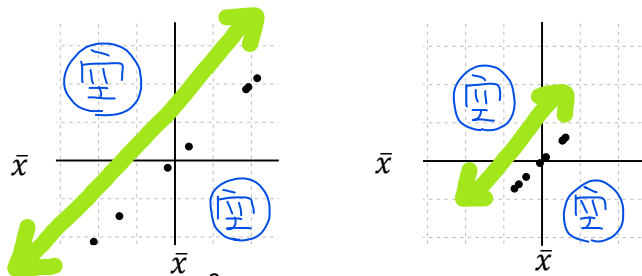
第3案  
都合のいい回避

1つのデータではなく  
データ全体を見よう



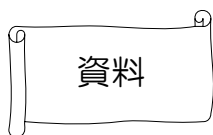
### 共分散の特殊なケースとしての分散

xとyという二変数の間ではなく、xとxという自分自身との間の共分散を考えてみたものが分散である。



$$\text{分散 } S_x^2 = S_{xx}$$

この値が大きいということは、その変数のばらつきが大きいということ。そこで、分散はばらつきの指標として使われる。



## 資料2-1 様々な「平均」

①算術平均 (arithmetic mean)、別名：相加平均

これは、もっとも一般的な「平均」で、すべての値を足し、総数で割ったもの。

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

②加重平均 (weighted average)

これは、各値にそれに対応する重みを付けたもの。

$$\bar{x} = \sum_{i=1}^n w_i x_i$$

算術平均は、次のように書けるので、加重平均の特殊な場合とみなすことができる。あるいは、加重平均は、重みが  $w_i = 1/n$  だという制約を取り払った、一般化された算術平均だとも言える。

$$\bar{x} = \sum_{i=1}^n \frac{1}{n} x_i$$

③調和平均 (harmonic average)

これは、逆数の平均。統計学では、分散の逆数（これを精度と呼ぶ）という概念がしばしば登場する。

$$\bar{x}_{\text{調}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$$

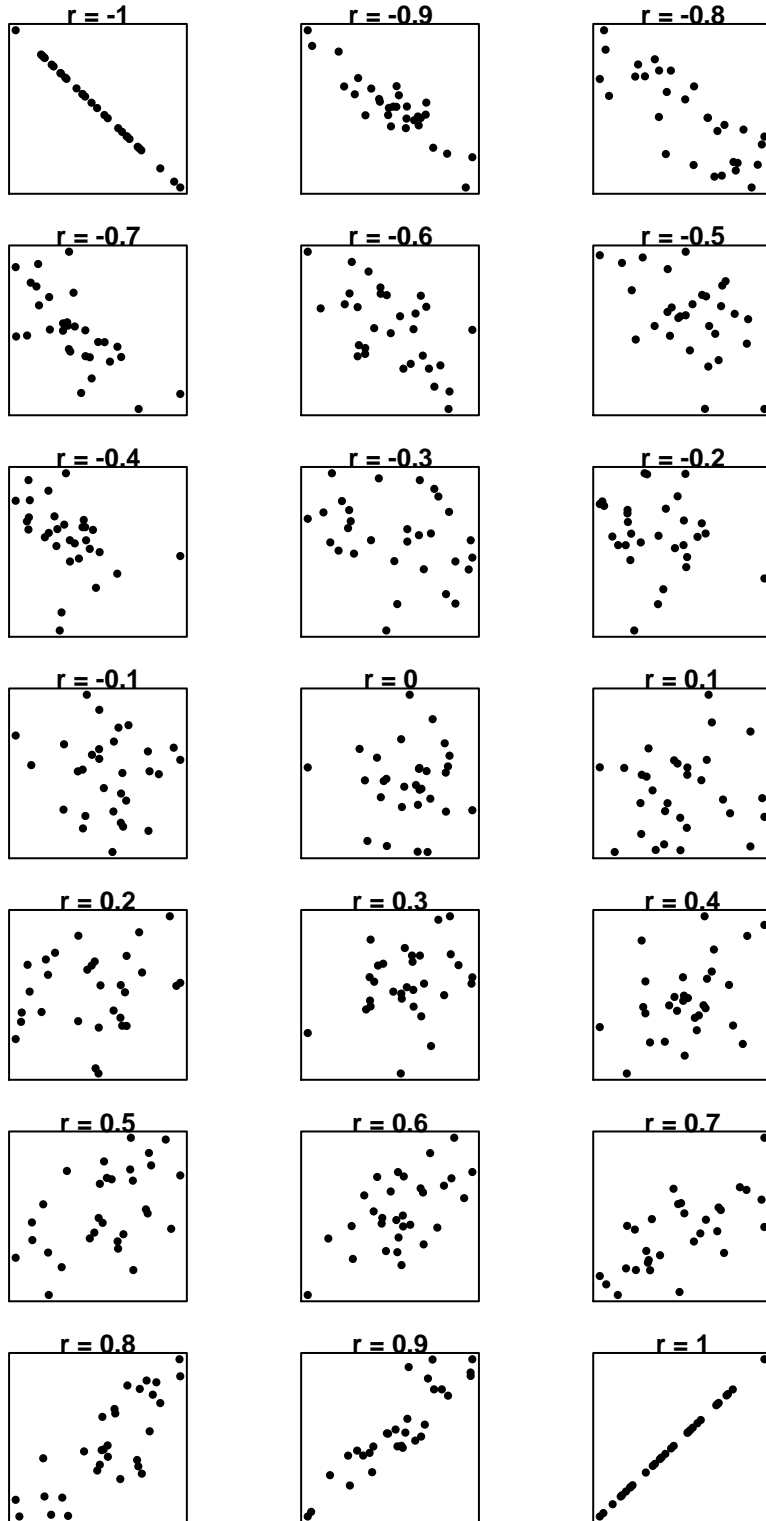
④幾何平均 (geometric average)、別名：相乗平均

これは、積に関する平均。つまり、一回当たり、平均してどれくらいの値をかけていたのかを表す。

$$\bar{x}_{\text{幾}} = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

$$\text{例} : (2 \times 4 \times 8)^{\frac{1}{3}} = (4 \times 4 \times 4)^{\frac{1}{3}} = 4$$

資料2-2 共分散／相関係数と散らばり



資料2-3 分散の分解

$$s_{x+y}^2 = s_x^2 + 2s_{xy} + s_y^2$$

【証明】

$$\begin{aligned} s_{x+y} &= \frac{1}{n} \sum_{i=1}^n \{(x_i + y_i) - (\bar{x} + \bar{y})\}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \{(x_i - \bar{x}) + (y_i - \bar{y})\}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \{(x_i - \bar{x})^2 + 2(x_i - \bar{x})(y_i - \bar{y}) + (y_i - \bar{y})^2\} \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{1}{n} \sum_{i=1}^n 2(x_i - \bar{x})(y_i - \bar{y}) + \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + 2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= s_x^2 + 2s_{xy} + s_y^2 \end{aligned}$$