

# 入門 統計学 第1章

## データの整理

『入門 統計学 第2版 一検定から多変量解析・実験  
計画法・ベイズ統計学まで』(オーム社)

※注: 本書を購入された方へのサービスですので, 教科書指定(参考図書は不可)していない授業での使用はお控えください。



# 統計学とは？

❁ 実験や調査によって観測されたデータ  
の特性を調べる学問

❁ 図表にまとめたり，統計量（平均や分散）を計算することで特性を把握・推定  
する



データの特性を調べる



# なぜ統計学を学ぶのか

- ❁ 大工さんが最初に学ぶのは、道具の使い方であるように、論文や報告書を書くための道具として統計学を学ぶ
- ❁ 実験結果の汎用性（偶然ではないこと）を示すことで、他人を納得させる



# 色々な場面で使われている統計学

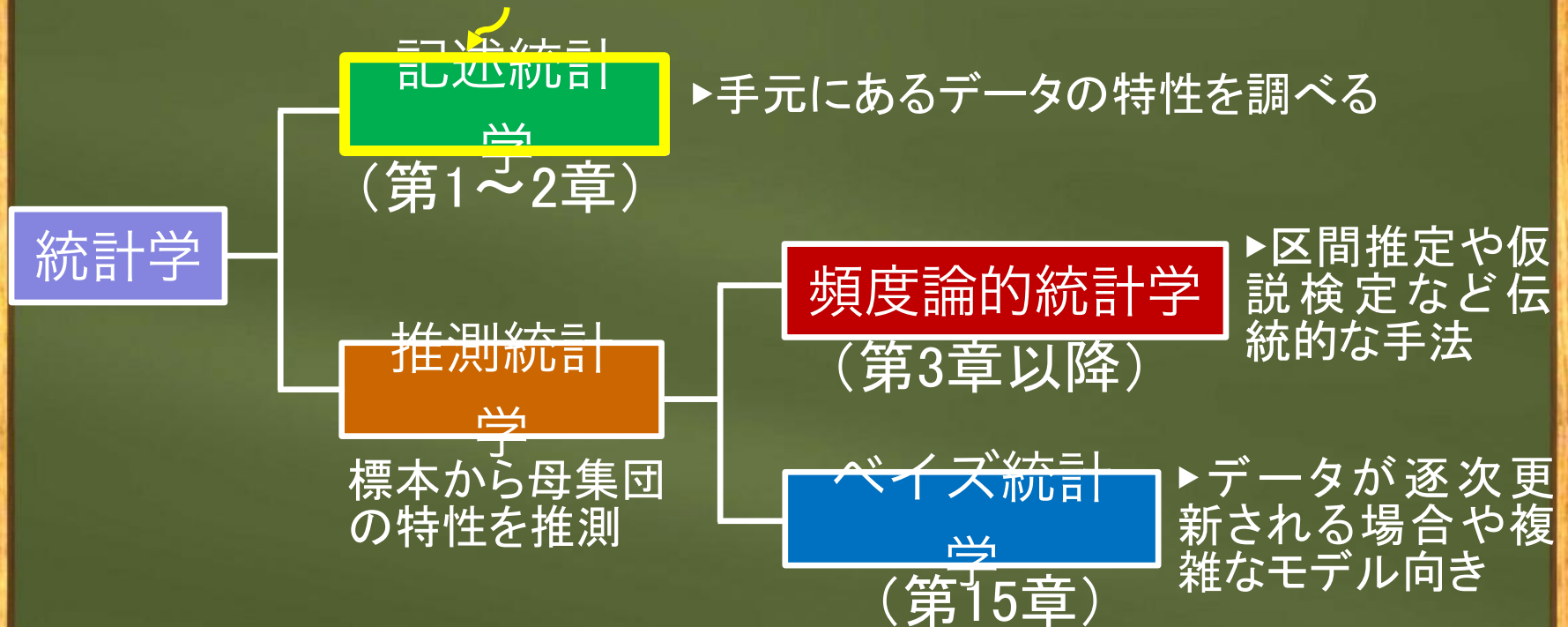
- ❁ 地震の予測
- ❁ 選挙速報の当確
- ❁ テレビの視聴率
- ❁ 新薬の有効性
- ❁ アパートの家賃
- ❁ 品質の管理
- ❁ ビッグデータの解析…，他にも沢山！



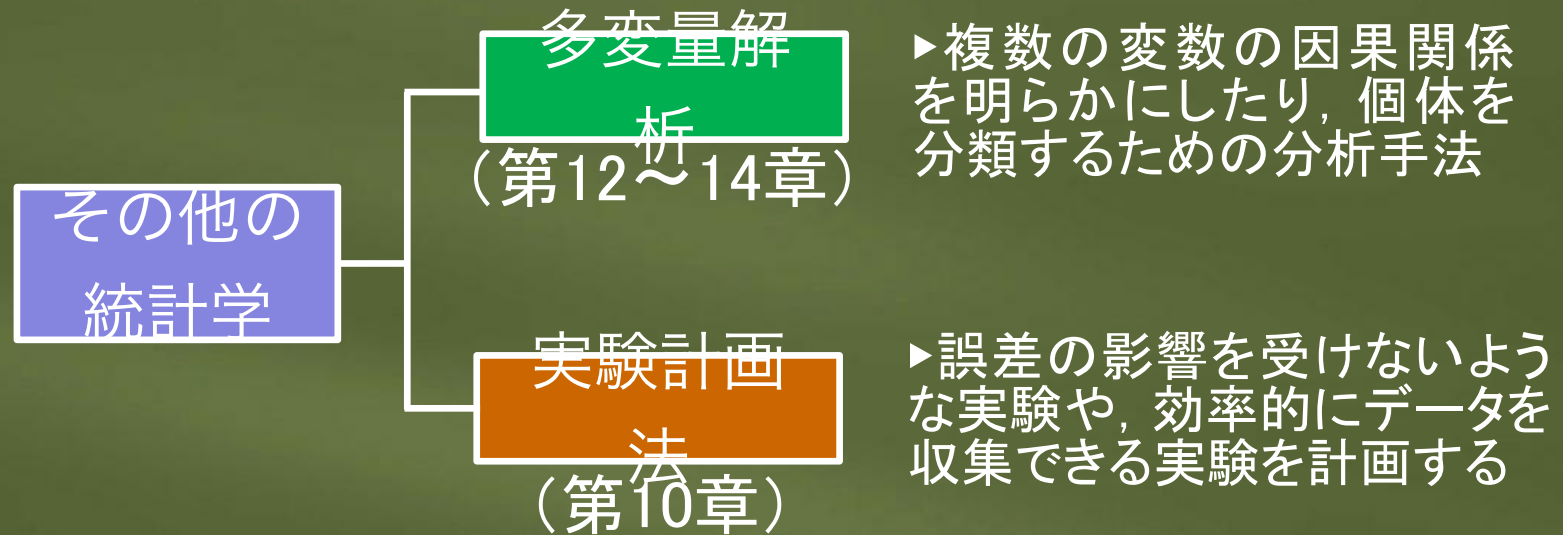


# 統計学の種類

推測統計学の基本でもあるので、これから学ぶ



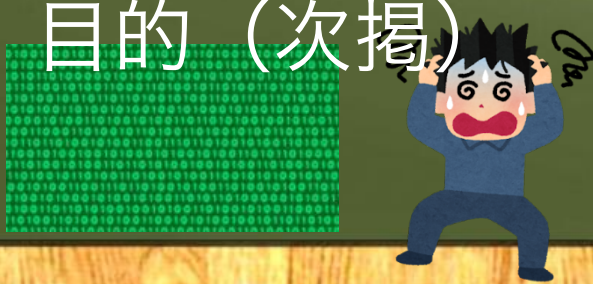
# その他の統計学



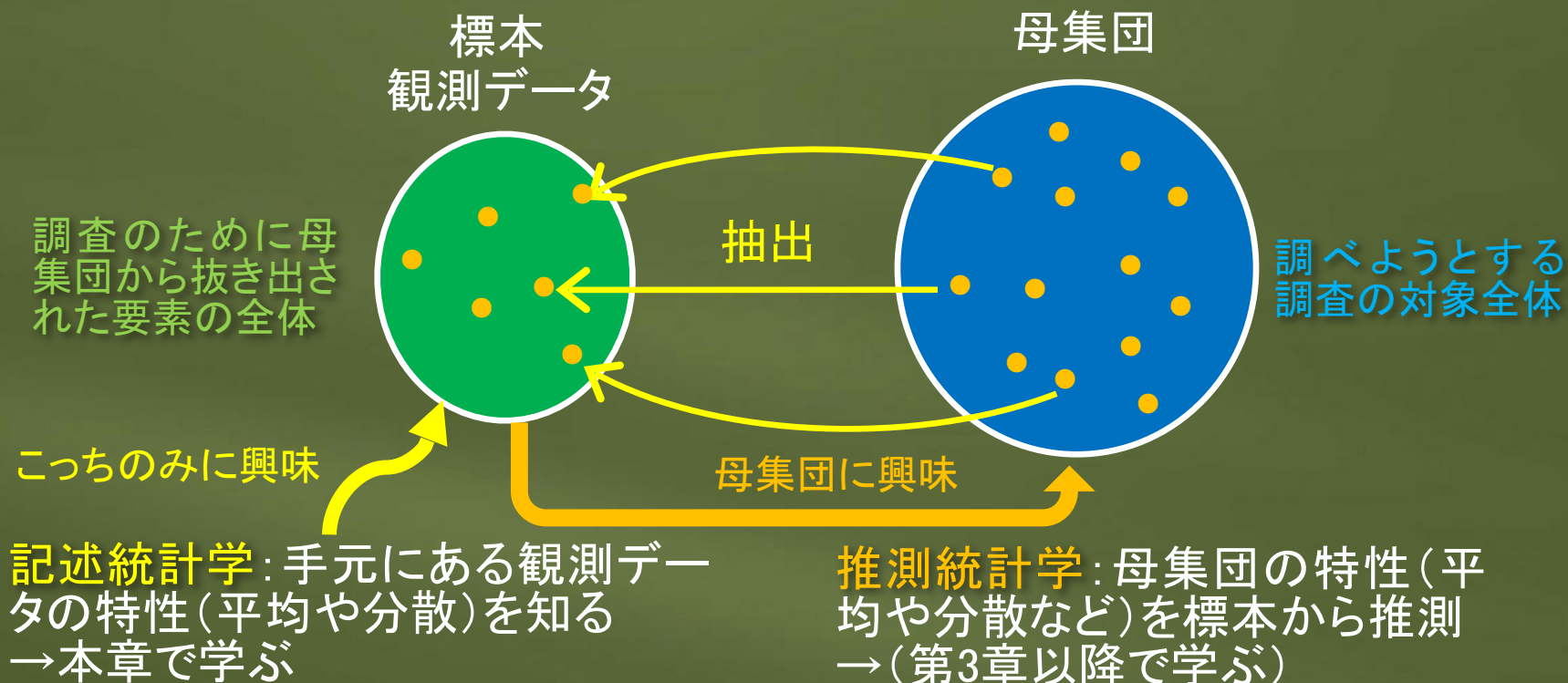


# 1.1 記述統計学

- 記述統計学：統計量（平均や分散）や図表に整理することで，観測したデータ（変数に値が入った状態）の特性を捉える統計学
- 推測統計学（第3章以降）との違い：記述統計学は手元のデータの特性を知ることが目的（次掲）



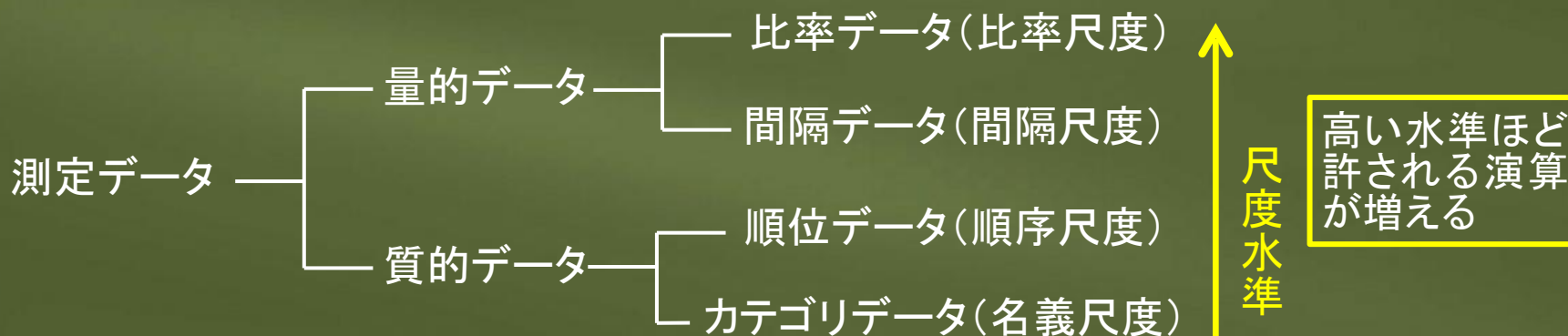
# 記述統計学 vs 推測統計学





# 尺度水準の問題

- ❖ データを測定（観測）する尺度によって、許される計算や分析が異なる（**尺度水準の問題**）
- ❖ 手元のデータは、どの尺度水準で測定されたものなのかを区別できなければならない



# 測定尺度① ー量的データー

❁ **比率データ**（質量，長さ，時間など）：  
絶対的なゼロを持っているため，値の間の比にも  
意味がある→**四則演算が可能**



❁ **間隔データ**（摂氏温度，知能指数など）：  
相対的なゼロしか持たないが，値の間隔は等しい  
→**足し算や引き算が可能**



❁ これら2種類のデータ（**量的データ**）が統計学の  
主な対象となる



# 測定尺度② — 質的データ —

🌀 **順位データ**（満足度，選好度など）：

値の大小関係にのみ意味がある

→ 中央値の計算は可能

🌀 **カテゴリデータ**（性別，職業など）：

値は内容を区別するだけ

→ カウントしたり最頻値のみ可能

🌀 これら2種類のデータ（**質的データ**）の分析には専用の手法が必要（第11章）

アンケート  
などで観測  
されるデー  
タね



# (まとめ) 測定尺度の一覧

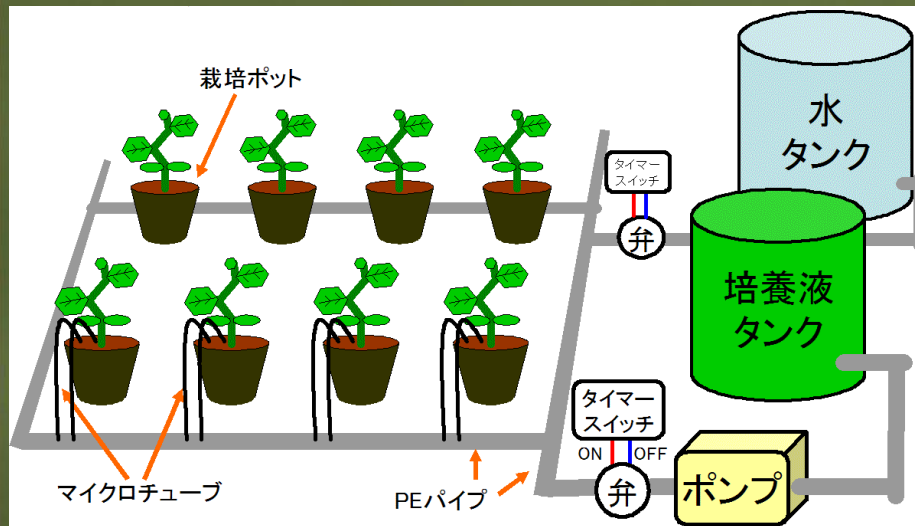
	データ名称	測定尺度	直接できる演算	主な代表値	事例
量的データ※	比率データ	比率尺度	+ - × ÷	幾何平均	質量
	間隔データ	間隔尺度	+ -	算術平均	摂氏温度
質的データ	順位データ	順序尺度	> =	中央値	満足度
	カテゴリデータ	名義尺度	度数カウント	最頻値	性別

※量的データには、質量のように値が連続している**連続型**と、金額のように飛び飛びになっている**離散型**がある(質的データは全て離散型)。



# 1.2 度数分布表とヒストグラム

キュウリの養液土耕栽培の実験



実験の目的:  
昼・夜どちらの時間帯に養液を与えた方が収量が多くなるのか(それとも変わらないのか)を検証する

表1.2 キュウリの収量(g)

ポット番号	栽培法A (昼)	栽培法B (夜)
1	3,063	3,157
2	2,275	2,707
3	2,089	3,270
4	2,855	3,181
5	2,836	3,633
6	3,219	3,404
7	2,817	2,210

観測

値を並べただけの状態から  
特性を捉えるのは難しい…  
そこで、表や図に整理！

11	2,140	2,938
12	1,757	3,286
13	2,499	2,920
14	2,093	3,332
15	2,073	3,478

# 度数分布表の作成手順

①データを昇順に並び替える

→Excelなどの表計算ソフトならば簡単

②階級の数と幅, 階級値を決める

→階級数は $\sqrt{n+1}$  か  $\log_2 n+1$ などを参考 (nは総度数)

③各階級の度数を求める

④適宜, 相対度数や累積相対度数を求める

→相対度数 =  $\frac{\text{その階級の度数}}{\text{総度数}}$



# 度数分布表 (キュウリの事例)

階級を代表する値で中央値を使うことが多い

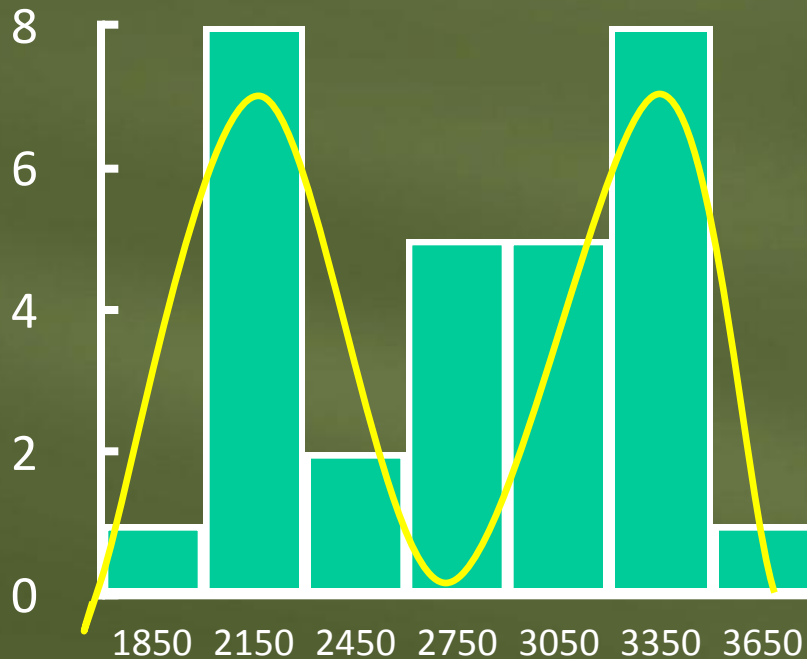
収量 (g)	階級値 (g)	度数 (ポット数)	相対度数 (%)	累積相対度数 (%)
1700以上～2000未満	1850	1	3.3	3.3
2000～2300	2150	8	26.7	30.0
2300～2600	2450	2	6.7	36.7
2600～2900	2750	5	16.7	53.3
2900～3200	3050	5	16.7	70.0
3200～3500	3350	8	26.7	96.7
3500～3800	3650	1	3.3	100.0

これら2つの階級で度数が多くなっていることが容易に確認できる

階級数は  $\sqrt{30+1}=6.5$ ,  $\log_2 30+1=5.9$  だが、切りの良い区間 (300g) とするため7階級と

# ヒストグラム

🔗 度数(ポット)数



度数分布表の階級を横軸X,  
度数を縦軸Yとして縦棒グラフを作成

二双峰性の分布になっていることから、養液を与える時間帯によって収量に差がありそう…

🔗 階級値



# Excel分析ツールによる 度数分布表とヒストグラムの作成

階級値(最大値)を適当な場所に入力しておく(最大値の数が階級数になる)

分析ツールを起動して[ヒストグラム]を選択

栽培法A	栽培法B	
3,063	3,157	2000
2,275	2,707	2300
2,089	3,270	2600
2,855	3,181	2900
2,836	3,633	3200
3,219	3,404	3500
2,817	2,219	3800
2,136	2,730	
2,540	3,408	
2,263	3,203	
2,140	2,938	
1,757	3,286	
2,499	2,920	
2,093	3,332	
2,073	3,478	

データ分析

分析ツール(A)

- 分散分析: 一元配置
- 分散分析: 繰り返しのある二元配置
- 分散分析: 繰り返しのない二元配置
- 相関
- 共分散
- 基本統計量
- 指数平滑
- F検定: 2標本を使った分散の検定
- フーリエ解析
- ヒストグラム**

ヒストグラム

入力元

入力範囲(I):

データ区間(E):

ラベル(L)

出力オプション

出力先(Q):

新規ワークシート(P):

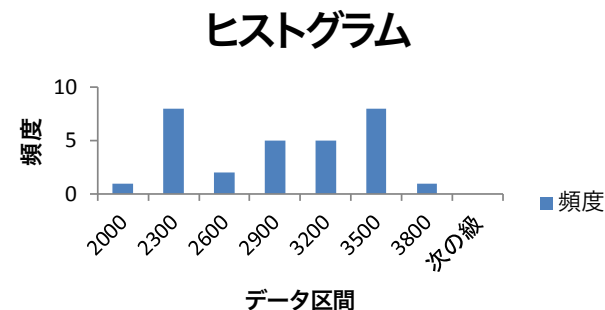
新規ブック(W)

パレート図(A)

累積度数分布の表示(M)

グラフ作成(C)

データ区間	頻度
2000	1
2300	8
2600	2
2900	5
3200	5
3500	8
3800	1
次の級	0



## 1.3 代表値① —平均—

🌸 代表値：データの特徴を表す統計量\*

※データに対して何らかの計算をして得られた値

🌸 平均や、バラツキを表す分散などがある

🌸 色々あるなかから、まずは平均をいくつか解説する（算術平均，加重平均，幾何平均，移動平均の順）



# 算術平均

## 🌻 小学校で学んだ誰もが知っている平均

平均を表す記号※  
(エックスバー)

算術平均

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

第 $n$ 番目のデータ値の意味  
 $n$ は添字(そえじ)・添数(てんすう)

$n$  ← データの数(標本サイズ) ← 標本数ではない

↓ 総和記号  $\Sigma$  (シグマ) を使う

$i=1 \sim n$ までの $x$ を足し合わせる  
( $\Sigma$ の上下記号は以降省略)

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

※第3章以降で学ぶ推測統計学では、 $\bar{x}$ は標本の平均という意味で、母集団の平均である $\mu$ と区別する。

Excel関数=AVERAGE

# 例題（キュウリ収量）

## 算術平均をExcel関数で求める

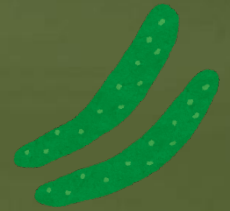


表1.2 キュウリの収量(g)

ポット番号	栽培法A (昼)	栽培法B (夜)
1	3,063	3,157
2	2,275	2,707
3	2,089	3,270
4	2,855	3,181
5	2,836	3,633
6	3,219	3,404
7	2,817	2,219
8	2,136	2,730
9	2,540	3,408
10	2,263	3,203
11	2,140	2,938
12	1,757	3,286
13	2,499	2,920
14	2,093	3,332
15	2,073	3,478

両栽培法合わせた算術平均  
=AVERAGE(B3:C17)=2784

栽培法A（昼）の算術平均  
=AVERAGE(B3:B17)=2444

栽培法B（夜）の算術平均  
=AVERAGE(C3:C17)=3124



# 加重平均

(群別標本サイズなどを重みとした平均)

観測された値 $x_i$ に重み $w_i$ をかけて平均を計算

加重平均  $\bar{x}_w = \frac{w_1x_1 + \dots + w_nx_n}{w_1 + \dots + w_n} = \frac{\sum w_i x_i}{\sum w_i}$

例：女子20名，男子10名のクラスで，女子の平均が65点，男子が80点のクラスの平均点は？（算術平均

は72.5点)  
$$\bar{x}_w = \frac{65 \times 20 + 80 \times 10}{20 + 10} = 70$$

(女子の人数が多かったため) 算術平均72.5点よりも下がる

# 幾何平均

(外れ値のあるデータや変化率の平均)

🔍 観測値をかけ合わせた積のn乗根

幾何平均  $\bar{x}_g = \sqrt[n]{x_1 \times \cdots \times x_n} = \sqrt[n]{\prod x_n} = \left(\prod x_n\right)^{\frac{1}{n}}$

総乗記号 (パ

特徴：外れ値に引っ張られないため、極端な値の出やすい細菌数の測定（の平均）によく使われ

Excel関数=GEOMEAN(数値)



# 幾何平均で変化率の平均を求める

❶ 掛け算の累乗根である幾何平均は、掛け算である「変化率の平均」にも適している

❷ 例（物価）：2年前から昨年にかけて物価が2倍になり、昨年から今年にかけては8倍になりました。この2年間の物価の対前年比の平均は何倍でしょうか？

算術平均： $(2+8) \div 2=5$  or 幾何平均： $\sqrt{2 \times 8}=4$

答え：2年前に100円のものが昨年は200円になり、今年は1600円（16倍）になっています。算術平均を2回（2年分）掛けると25倍、幾何平均を2回だと16倍なので、後者の方が妥当でしょう。

# 移動平均



(変動を抑えて傾向を見る)

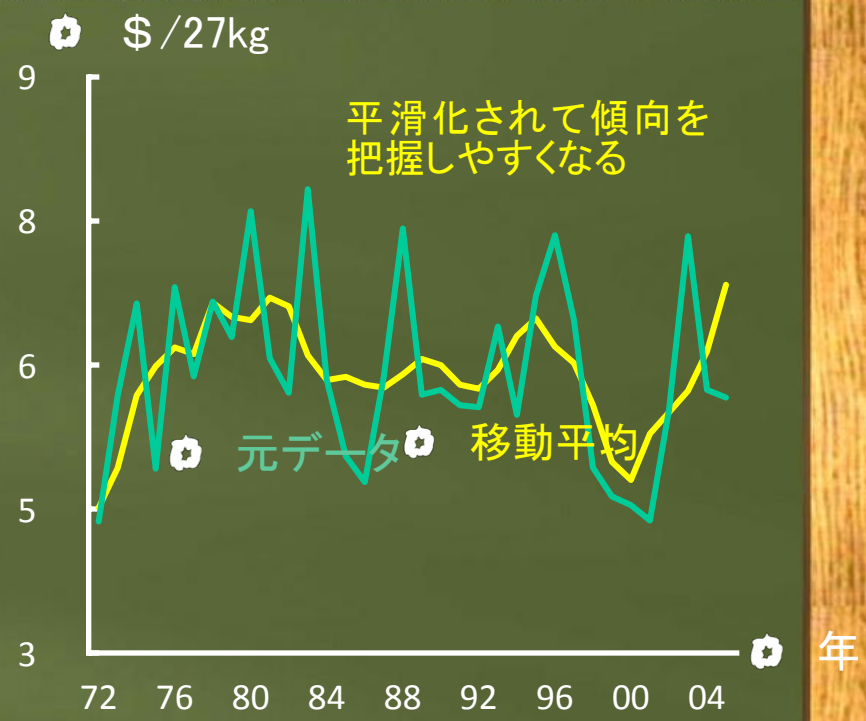
	大豆価格 (シカゴ)	移動平均 (5項)
1970	2.85	
1971	3.03	4.51
1972	4.37	4.93
1973	5.68	5.68
1974	6.64	5.99
1975	4.92	...
...	...	5.73
2003	7.34	6.14
2004	5.74	6.83
2005	5.66	
2006	6.43	
2007	9.00	

1970~1974の5年間※の平均

1年ずらした5年間の平均

同様に、1年ずらした平均

最初と最後の2年ずつ失う



※何項(年)の平均を取るかは自分で決める

Excel分析ツールにも「移動平均」が搭載



## 1.4 代表値② —バラツキの指標—

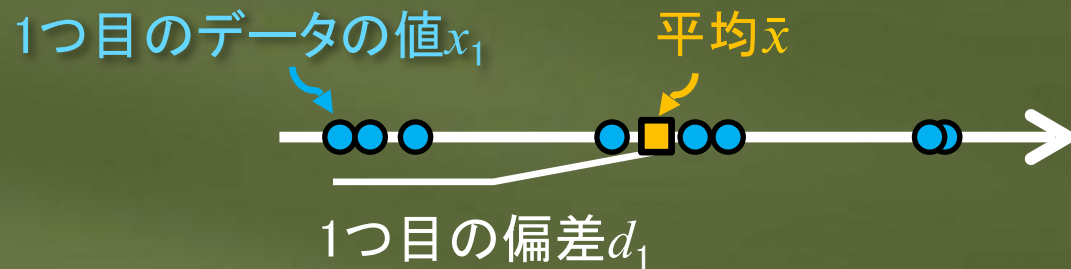
- ❁ データの特性として平均だけでは情報不足
- ❁ どのぐらいバラついているのかも知りたい
- ❁ バラツキの指標（統計量）として,  
偏差 → 偏差平方和 → 分散 → 標準偏差  
→ 変動係数, の5つを順に解説

# 偏差

(もっとも基本的なバラツキの指標)

🗄 各データ値と平均との差 (個別データが平均からの偏り度)

偏差  $d_i = x_i - \bar{x}$



偏差の欠点:

- ① データの数だけ示さなければならない
- ② 統計量として1つの値にまとめようと足し合わせると、±値が混在しているため、相殺されて0になってしまう



# 偏差平方和と分散

- ❖ 偏差を2乗して足し合わせれば0にならない

**偏差平方和**(変動)  $S = \sum (x_i - \bar{x})^2$

Excel関数=DEVSQ

**欠点**：標本サイズnが大きいと巨大な値になってしまう

- う
- ❖ 標本サイズで割れば解決（サイズに応じて小さくなる）

**分散**※  $s^2 = \frac{\sum (x_i - \bar{x})^2}{n}$

Excel関数=VAR.P

**欠点**：単位が2乗のままなので，値はまだ大きい

※推測統計学では，母集団の分散は $\sigma^2$ で表して区別します。

# 標準偏差

- ❖ 分散の平方根を取ることで、元の単位に戻り、大きさもちょうど良くなる

標準偏差 (SD)  $s = \sqrt{\text{分散}} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$  Excel関数=STDEV.P

欠点：値の大きさが極端に異なる集団間や、単位が異なる集団間のバラツキは比較できない



平均値の大きな集団の標準偏差が大きくなってしまふ



# 変動係数

❖ 標準偏差を平均で割ることで、集団間の平均を揃えて、単位を取り去る（無名数にする）

$$\text{変動係数 } CV = \frac{\text{標準偏差 } s}{\text{平均 } \bar{x}} = \frac{\sqrt{\sum (x_i - \bar{x})^2 / n}}{\bar{x}}$$

**注意**：5種類のバラツキの統計量は、それぞれ有する情報量が異なるので、分析に適したものを使い分ける  
（例：偏差ならば分布形もわかるが、変動係数は単位さえわからない）

# 1.5 質的データの代表値

- 質的データは平均 $\bar{x}$ の計算が許されない（分散も不可）
- 平均に代えて、カテゴリーカルデータでは最頻値を、順  
位データでは最頻値や中央値を計算する

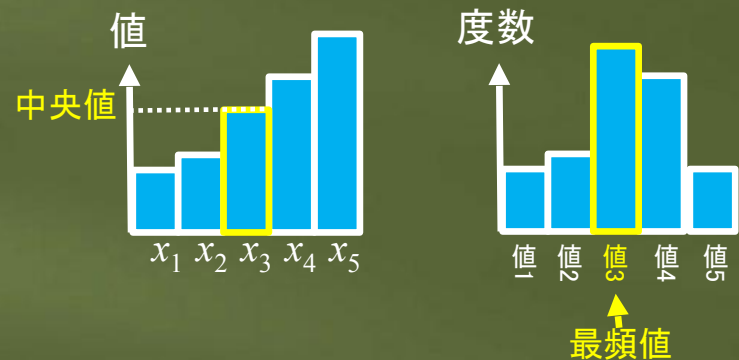
- 中央値**：データを大きさ順に並べたとき中央にく

Excel関数=MEDIAN

- 最頻値**：最も多く現れる（最大度数を持つ）

Excel関数=MODE

- バラツキ**は図や表に表して視覚的に確認する

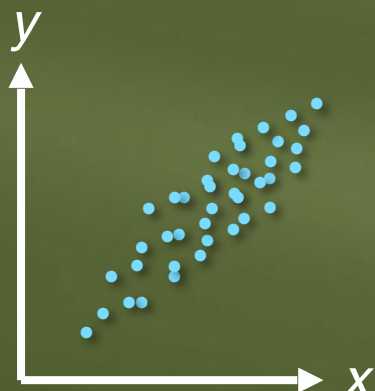




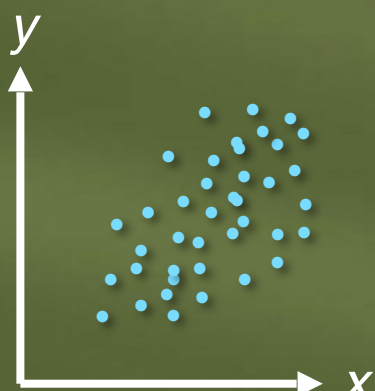
# 1.6 相関係数と共分散

## — 2つの変数の関係 —

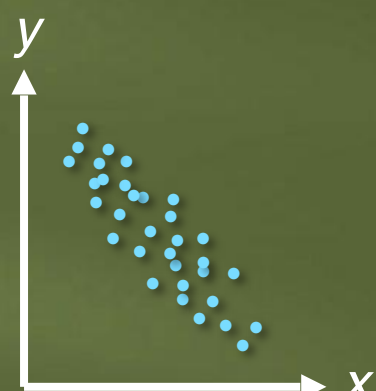
- ❶ 2変数間で、片方が大きくなるに従ってもう片方も大きくなる場合、「正の相関関係がある」という（小さくなる場合は負の相関関係）
- ❷ 相関関係の強さを示す統計量が相関係数と共分散



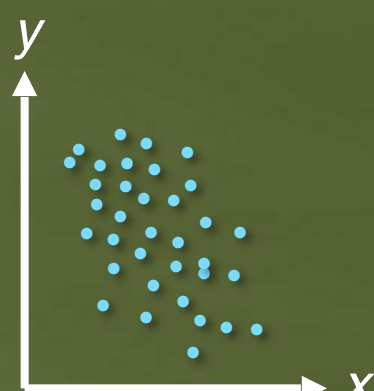
強い正の相関関係



弱い正の相関関係



強い負の相関関係



弱い負の相関関係

# 共分散

🔗 2変数の偏差の積の平均で相関関係を捉える

偏差の積の総和は、正の相関があるときには正となる（負の相関があるときには負となる）

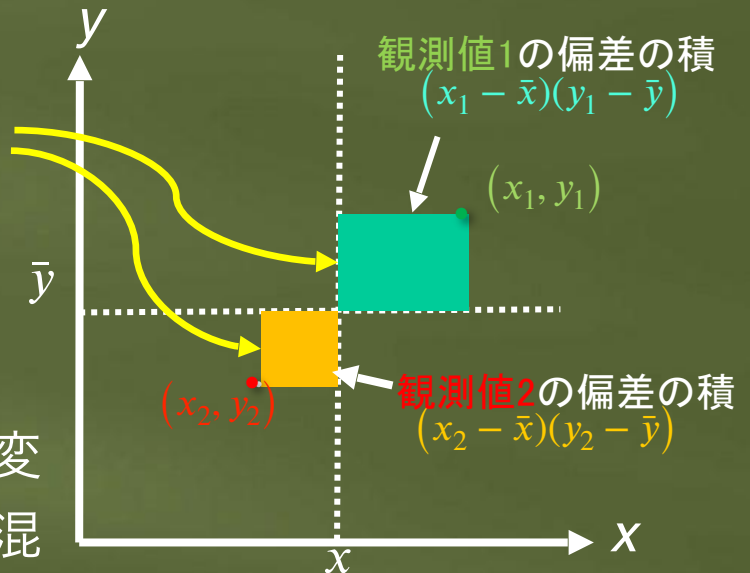
$$\sum (x_i - \bar{x})(y_i - \bar{y})$$

標本サイズに左右されないように  $n$  で

割る

$$\text{共分散 } s_{xy} = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

欠点：バラツキによって値が大きく変わってしまう（異なる単位も混



Excel関数⇒ COVARIANCE.P



# 相関係数

❖ どちらの偏差も標準偏差 $s$ で割ってバラツキをそろえる

$$\text{相関係数 } r_{xy} = \frac{1}{n} \sum \frac{(x_i - \bar{x})}{s_x} \cdot \frac{(y_i - \bar{y})}{s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

相関係数の性質：

- ① -1から1までの値を取る（0は無相関）が、いくつ以上が強い相関という基準はない
- ② あくまで2変数の直線的な関係を捉えるだけ

（右図のような関係があっても捉えられない）  
Excel関数=CORREL  
い→



以上で第1章は終了です。