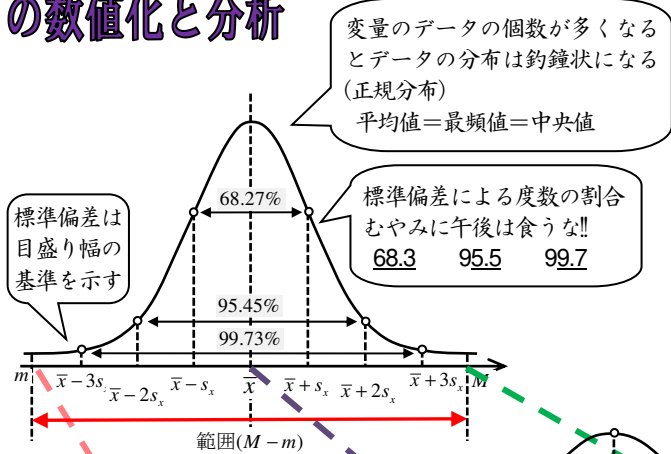


データの散らばりの数値化と分析

変量 (X)
計測対象の特性(プロパティ)
・数学の試験 …… 点数
・身体測定 …… 身長, 体重
・さいころ …… 目の値
・天気 …… 気温, 降水量

データ (data)
変量の測定値の集まり
 $x_1, x_2, x_3, \dots, x_n$
(データの大きさ n)



変量 X のデータ
 $x: x_1, x_2, x_3, \dots, x_n$
の散らばりの数値化ステップ

Step1
データの平均値 \bar{x} を求める
$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

Step2
平均からの散らばり(偏差)
 $x_1 - \bar{x}, x_2 - \bar{x}, x_3 - \bar{x}, \dots, x_n - \bar{x}$

偏差の和は0だから、和で散らばりは表せない

Step3
散らばりを偏差の大きさの和で数値化
○絶対値で散らばりを表す
 $|x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_n - \bar{x}|$
○2乗して散らばりを表す
 $(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2$

2つの変量のデータの個数が違うと、散らばりの比較はできない。

変量の変換

変量 X のデータ x を a 倍し b を加えるとデータの散らばりはどう変化するだろう?

○ a 倍する ($a > 0$)
平均と標準偏差は a 倍になる

○ b を加える
平均は b 増える。散らばりは度数分布の階級がズレるだけなので変わらない。

$x \rightarrow y = ax + b$ とすると
平均 $\bar{y} = a\bar{x} + b$
標準偏差 $s_y = |a|s_x$
 $\Rightarrow s_y^2 = a^2s_x^2$

Step4
データの個数で割り、散らばりを揃える(平均化)
○平均偏差(Mean Deviation)
$$MD = \frac{|x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_n - \bar{x}|}{n}$$

○分散(Variance)
$$s_x^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$$

分散の単位はデータの単位の2乗になってしまう
例えば、長さ(cm)の分散は面積(cm^2)

Step5
データと分散の散らばりの単位を揃える
○標準偏差(Standard Deviation)
$$s_x = \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}}$$

○偏差値(Standard Score)
$$ss = \frac{x - \bar{x}}{s_x} \times 10 + 50$$

目盛の変換
 $s \Rightarrow 1 \Rightarrow 10$
平均の変換
 $\bar{x} \Rightarrow x - \bar{x}(0) \Rightarrow 50$

仮説検定
有意水準(危険率)付きの背理法

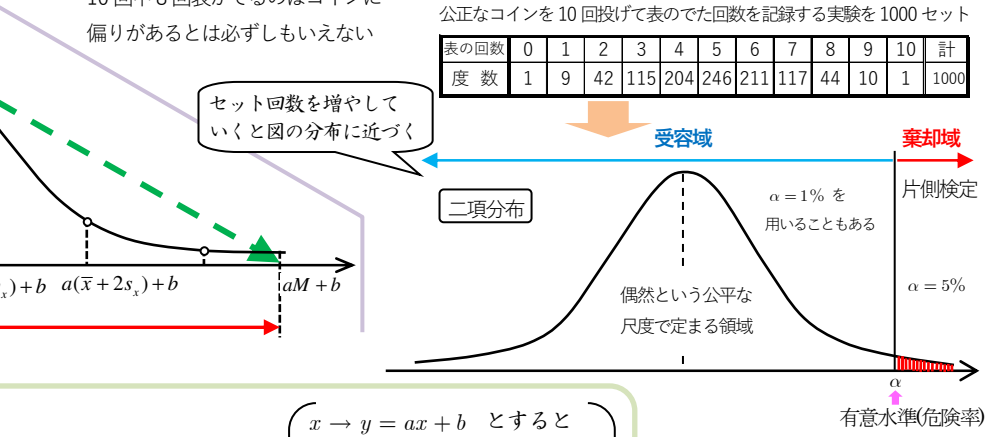
対立仮説 H_1 vs 帰無仮説 H_0

主張が正しきことはどのように判定する?
主張は正しくない
公平な実験を繰り返し統計量(度数)を調べる
有意水準による判定

主張は正しくないという判定は積極的に採択されるわけではない
帰無仮説 H_0 は棄却される(無に帰する) → 棄却域
帰無仮説 H_0 は採択される → 受容域

コインを投げて10回中9回表がでるのはコインに偏りがある?
9回以上の相対度数は $\frac{10+1}{1000} = 0.011 < \alpha$

有意水準を $\alpha = 5\%$ ($= 0.05$) とする。
10回中9回表がでるのはコインに偏りがある



仮平均 $x_0, u = \frac{x - x_0}{c}$ とする
 $u = \frac{x - x_0}{c}, s_u^2 = \frac{1}{c^2}s_x^2$

Fuminori Nakamura

偏差を計算しないで標準偏差(分散)は求められる

○標準偏差の簡便式
$$s_x = \sqrt{x^2 - (\bar{x})^2}$$