

2017年度日本疫学会スライドコンテスト受賞作品

# 基本的な統計 —データの分布とばらつき—

藤田保健衛生大学医学部

公衆衛生学

柿崎 真沙子

# 基本的な統計

- ① 母集団と標本
- ② 誤差・ばらつきを表す要約値
- ③ データの分布
- ④ 標準化と偏差値

# 母集団と標本

- **母集団**：調査対象とする（推測結果を適用したい）個体の集まり、集団
  - 母集団は一般的に全数調査が困難
    - 例：日本人の平均身長を知りたい！
- **標本**：無作為に抽出された個体
  - 例：電話帳から100名無作為に抽出
- A県在住の地域住民（**=母集団**）の生活習慣を調査したいのに、特定の疾患を持つ患者だけを調査しても偏りが生じてしまう（**=誤差**）
- 母集団から偏りなく標本を集めることが重要

# 例) 2群の平均を比較する

- 某アイドルグループ全体を母集団とした場合、チームAは標本として適切か？
- チームAとチームBの平均年齢に違いは有るか？

**どうやって検討したらいいだろうか？**  
(その前段階として分布とばらつきを検討)

# まずは誤差とばらつき

- 誤差・ばらつき
- 誤差・ばらつきを表す要約値
  - 平均
  - 偏差
  - 偏差平方和
  - (不偏) 分散
  - 標準偏差(S)
  - 標準誤差

# 誤差・ばらつき

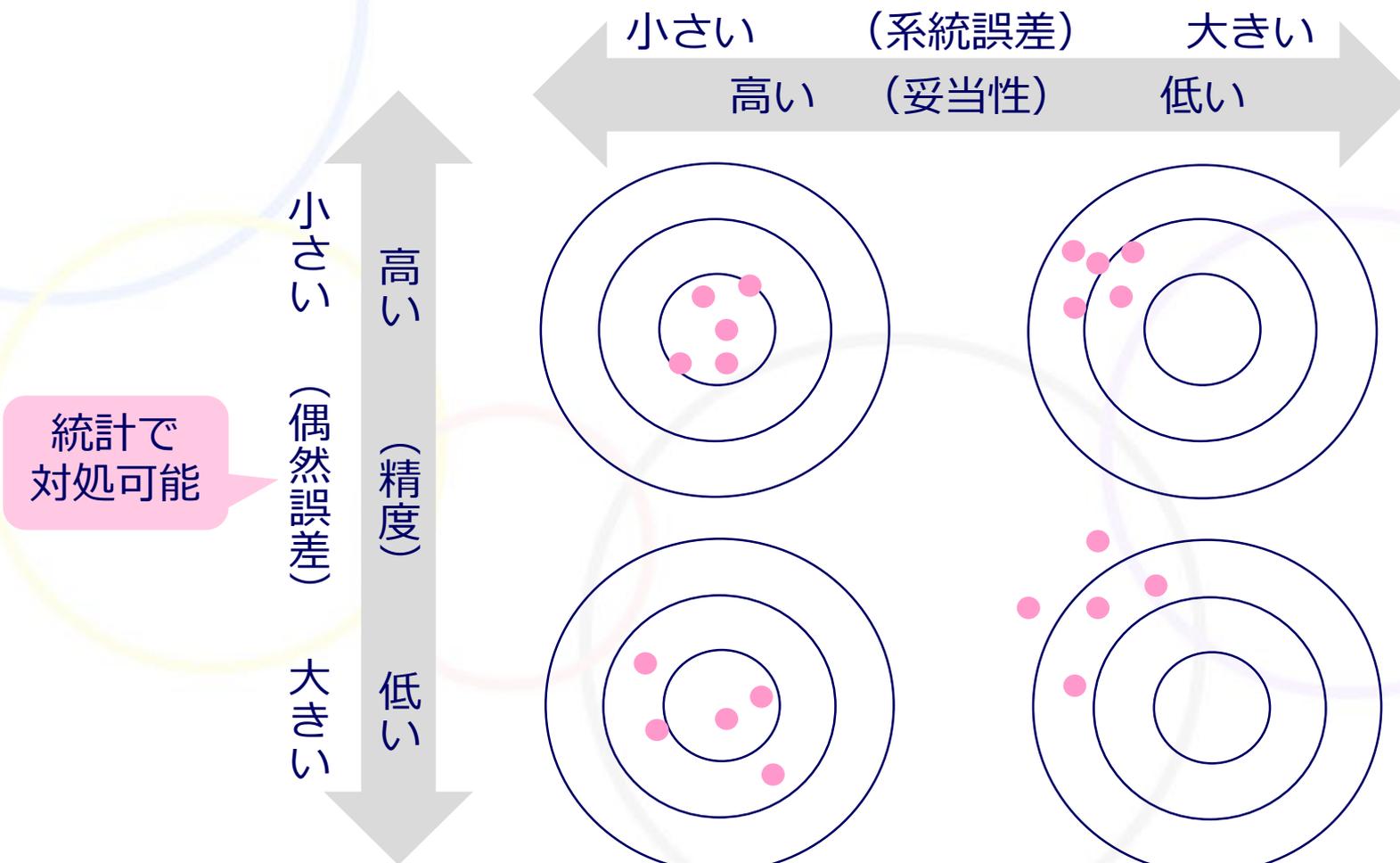
## • 偶然誤差 (Random error)

- 偶然によってのみ生ずるばらつき
- 真の値を中心にばらつきが分布
- 偶然を支配する法則 (統計学) により対処可能

## • 系統誤差 (Systematic error)

- 中心が真の値から偏っている
- バイアス (Bias) とも言う
- 統計学では対処不能

# 誤差・ばらつきが多い少ない



# どうやって 誤差・ばらつきを表すか？

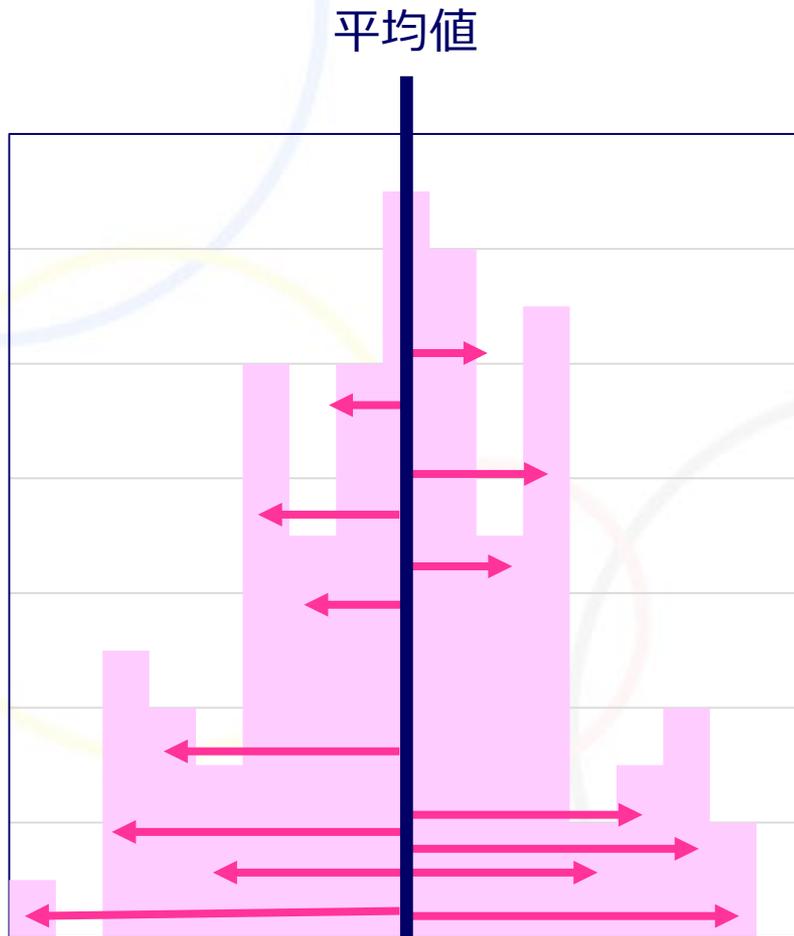
- 平均
- 分散 (Variance)
- 標準偏差 (Standard deviation : SD)
- 標準誤差 (Standard error : SE)
- 分位点
- パーセンタイル
- 範囲 (Range)

# ばらつきを表す要約値

N個のデータをそれぞれ、 $X_1, X_2, X_3 \dots X_n$ とする

- **平均 ( $\bar{X}$ )** :  $(X_1 + X_2 + X_3 + \dots + X_n) / n$
- **偏差** : (データ - 平均)
  - $(X_1 - \bar{X}), (X_2 - \bar{X}), (X_3 - \bar{X}), \dots (X_n - \bar{X})$
- **偏差平方和** : 偏差の二乗を全てのデータで合計
  - $(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + (X_3 - \bar{X})^2 + \dots (X_n - \bar{X})^2$

# 分散



- 平均から個々のデータがどのくらい離れているか？
- 「平均からの離れ具合の平均値」が分散
- 偏差平方和を $n-1$ で割ったもの

# 標準偏差

- **分散の平方根**
- 標準偏差はデータそのものの分布の広がり幅 (ばらつき) をみる尺度
  - 散らばり具合を見る
- 平均値と標準偏差：この2つがわかれば、データの分布 がある程度明らかに

# 標準誤差

- 真の平均値と標本の平均値のズレが標準誤差
- 標準誤差の値が大きいほど推定精度が低い
- 日本人全員の平均身長を知りたい！
  - 真の平均値
- 電話帳からランダムに100名選択
  - 標本の平均値



ズレ

# 例) チームAとチームBの年齢

チームA	チームB
20	24
25	28
23	24
20	22
28	21
23	25
22	26
24	27
25	27
18	23
27	27
20	22
23	28
29	23
17	25

# では計算してみよう！

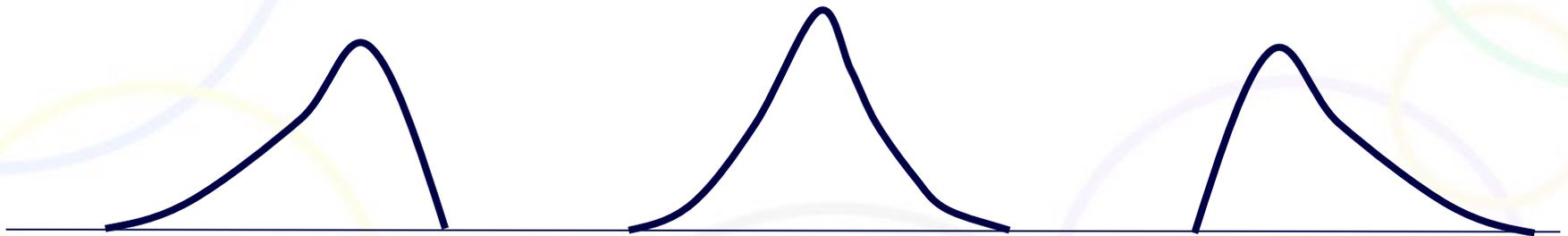
		チームA	チームB
平均年齢	すべて足してNで割る		
偏差平方和	(データー平均) <sup>2</sup> を全てのデータで合計		
(不偏) 分散	偏差平方和をn-1で割ったもの		
標準偏差	分散の平方根		
標準誤差	標準偏差を $\sqrt{n}$ でわる		

# では計算してみよう！

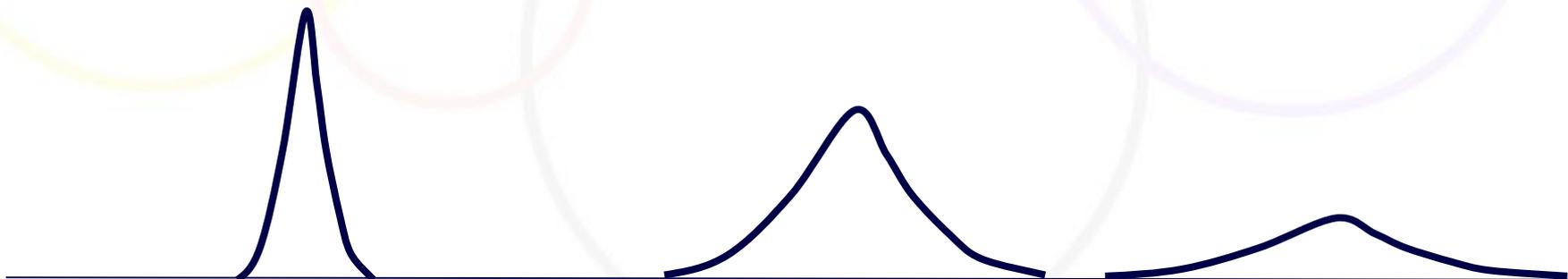
		チームA	チームB
平均年齢	すべて足してNで割る	$344/15 = 22.93$	$372/15 = 24.80$
偏差平方和	(データ-平均) <sup>2</sup> を全てのデータで合計	$(20-22.93)^2 + \dots = 174.933\dots$	$(24-24)^2 + \dots = 74.4$
(不偏) 分散	偏差平方和をn-1で割ったもの	$174.933/14 = 12.495$	$74.4/14 = 5.314$
標準偏差	分散の平方根	$\sqrt{12.495} \dots = 3.535$	$\sqrt{5.31} \dots = 2.305$
標準誤差	標準偏差を $\sqrt{n}$ でわる	$\sqrt{3.534}/\sqrt{15} = 0.9127$	$\sqrt{2.227}/\sqrt{15} = 0.5952$

# データの分布を表す指標

歪度 (ひずみ、skewness) : 分布の左右対称度



尖度 (とがり、kurtosis): 分布の裾の長さ



✓ 「正規分布」と言っても様々な分布の形がある

分布の形が異なると比較がしにくい



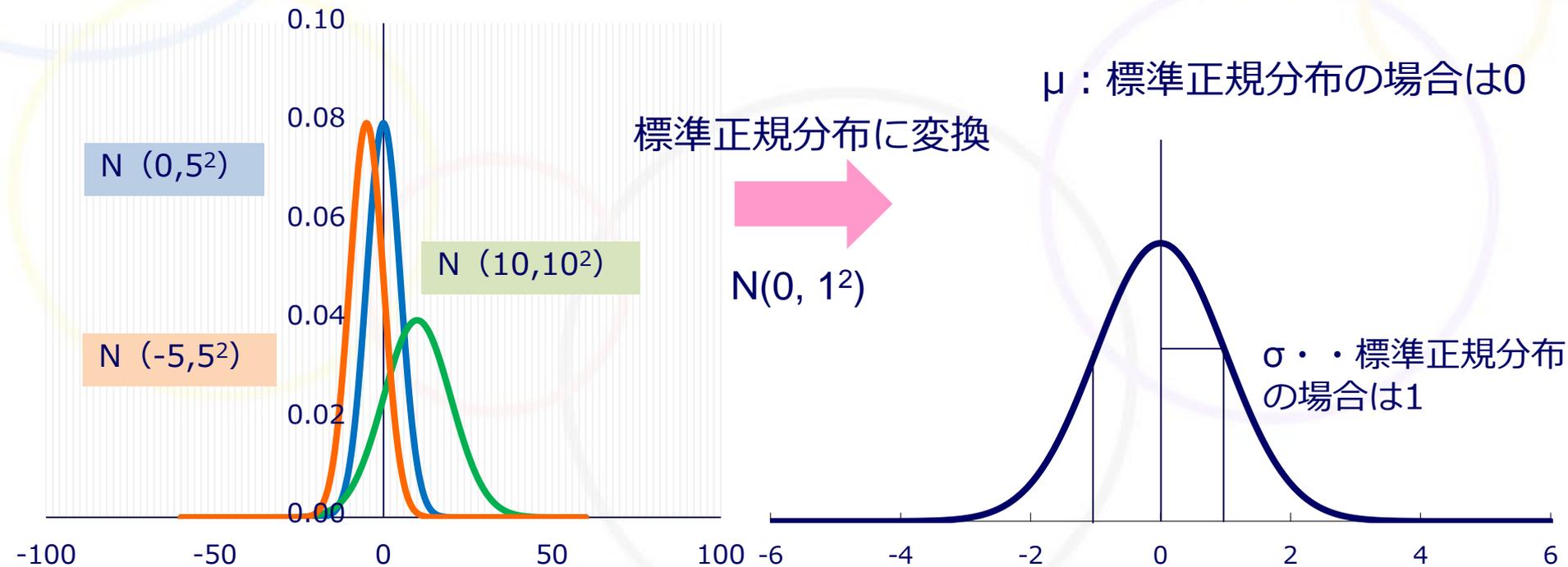
形をそろえて比較しやすくする



**標準化（z変換）**

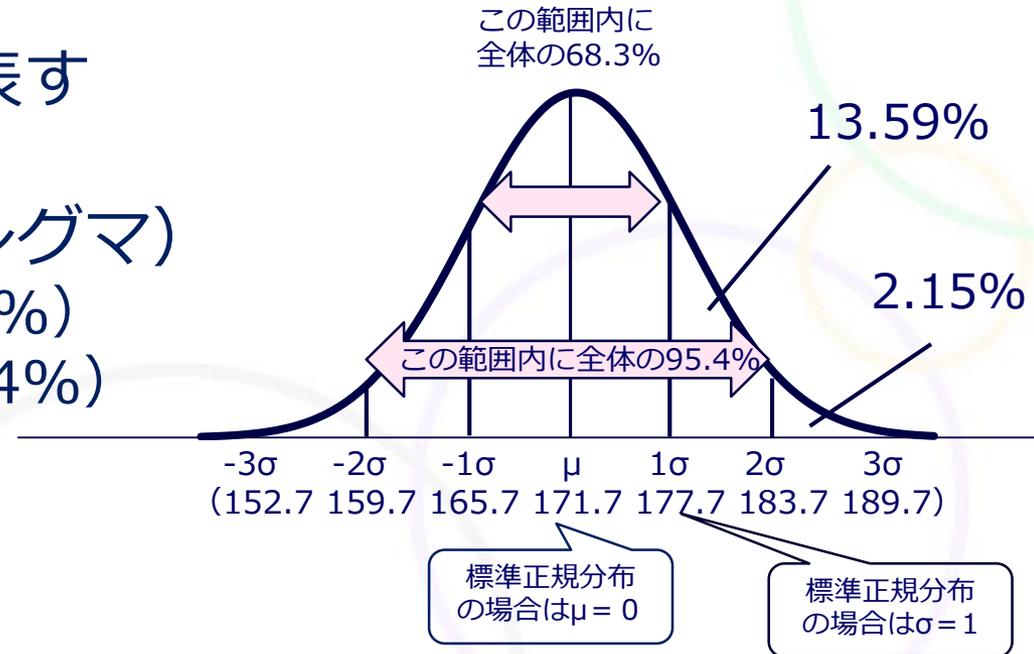
# 標本分布を標準化（Z変換）

- $N(\mu, \sigma^2)$  を  $N(0, 1^2)$  に変換できる
- 全てのデータを標準化した時、変換後の値を標準得点（z 値）と呼ぶ
- Z値 = 「どのくらい平均から離れているか」を示す。



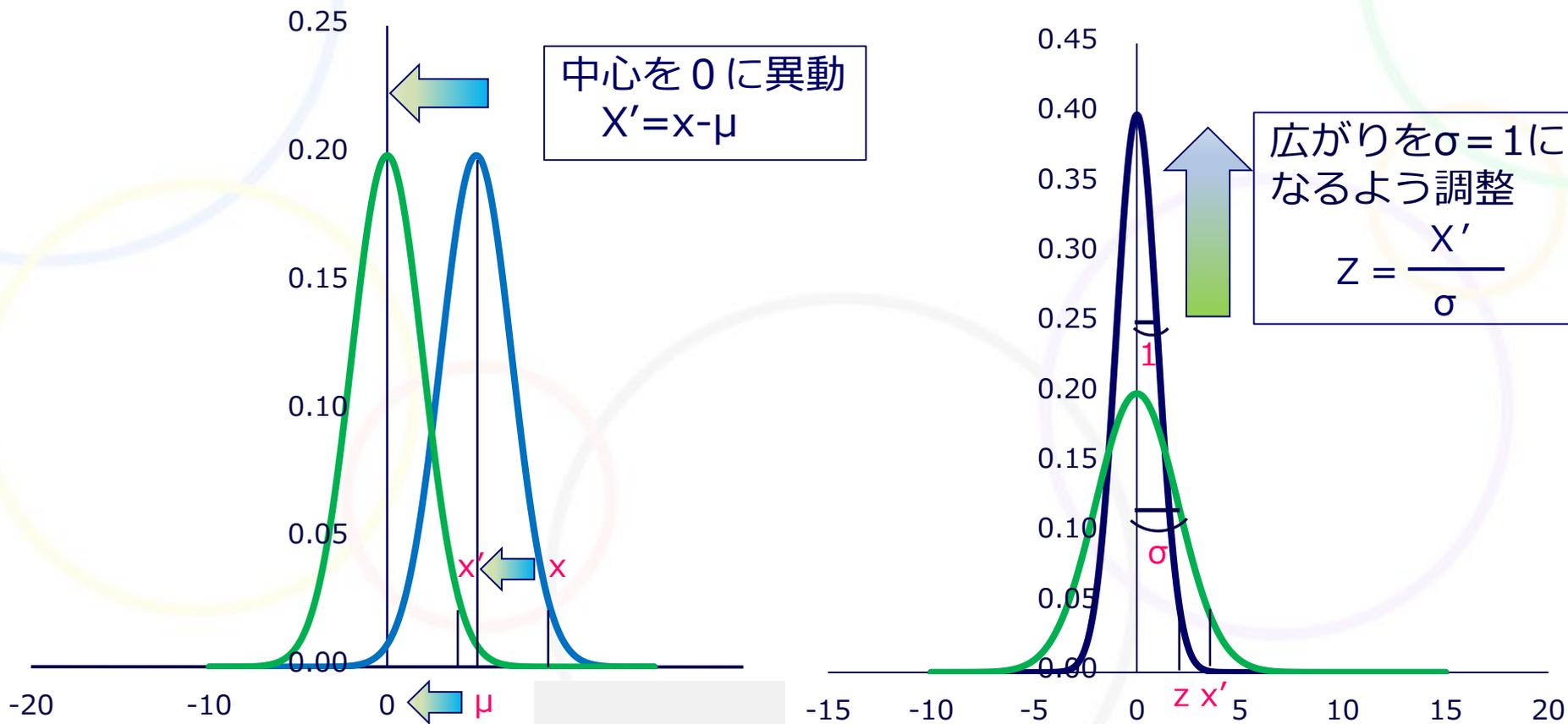
# (おさらい) 正規分布とは

- 正規分布  $N(\mu, \sigma^2)$  と表す
- 平均値:  $\mu$  (ミュー)
- 標準偏差 (SD):  $\sigma$  (シグマ)
  - $\pm 1SD$  に約  $2/3$  (68.3%)
  - $\pm 2SD$  に約 95% (95.4%)
  - $\pm 1.96SD$  に 95%
- 標準正規分布  $N(1, 0)$



- 例: 平成27年国民健康・栄養調査より26-29歳の男性の身長は平均171.7cm、標準偏差 ( $\sigma$ ) 6.0cmである。この度数分布が正規分布に従うとすると図のように示される。

# 標準化の手順 (イメージ)



# 標準化の公式

- すべてのデータを以下の公式を使って変換

$$Z = \frac{X - \mu}{\sigma}$$

①平均をひいて

②標準偏差で割る

- ① それぞれのデータから平均を引く
  - ② それを標準偏差で割る
- Zスコアはある値 $x_i$ が分布の中でどの辺りに位置するかを平均値0、標準偏差1の標準正規分布に置き換え、表したものの

# なぜ標準化？

- 定期テスト（英語90点、国語80点）
  - 英語の方がよくできたのか？？
- 平均点が異なるため一概に比較できない
- 標準化を行い（いわゆる）偏差値を出す
- 偏差値 = 平均値50、標準偏差10に標準化

$$\begin{aligned}\text{偏差値}T &= \frac{X - \mu}{\sigma} \times 10 + 50 \\ &= z \times 10 + 50\end{aligned}$$

# 偏差値を計算してみよう！

科目	平均値	標準偏差	得点	偏差値
英語	60	12	70	
国語	65	20	65	
数学	50	18	80	
理科	55	10	80	
社会	52	15	50	
5教科	282	75	345	

# 偏差値を計算してみよう！

科目	平均値	標準偏差	得点	偏差値
英語	60	12	70	58.3
国語	65	20	65	50.0
数学	50	18	80	66.7
理科	55	10	80	75.0
社会	52	15	50	48.7
5教科	282	75	345	58.4