

# 健康統計学 第10回

今回は、前回は引き続き、確率分布（テキスト53～63ページ）について学習します。とくに、正規分布（標準正規分布）、カイ二乗分布、t分布、F分布について説明します。

また、母集団統計値の推定（テキスト65～76ページ）についても学習します。

## テキスト

- 『やさしい保健統計学 改訂第4版』 縣 俊彦著(南江堂)

## 今回の内容

1. [正規分布と関係する分布](#)
  - a. [補足: 標準化を用いた確率の計算](#)
2. [母集団と標本](#)
3. [点推定と区間推定](#)

# 正規分布と関連する分布

自然界や一般に観察できる多くのものについて、その分布は、平均値を中心に左右対称の釣鐘状の分布になっていることがあります。

- 生物現象、毎年の雨量など
- 身長や体重、標準的なテストの成績など

## 正規分布 (normal distribution)

### 正規分布とは

- 平均値を中心とした左右対称の釣鐘状になる分布を、「正規分布」という

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- $\pi$  は円周率で、 $\pi = 3.14159\dots$
  - $e$  は自然対数の底で、 $e = 2.71828\dots$
  - $\mu$  : 平均値
  - $\sigma$  : 標準偏差
  - $\sigma^2$  : 分散
- 確率密度関数が上の式になるとき、『 $x$  は、平均が  $\mu$  で標準偏差が  $\sigma$  の正規分布にしたがう』といい、 $N(\mu, \sigma^2)$  と表す
    - つまり、平均値と分散(標準偏差)から分布が決まる
  - さまざまな推定や検定に使われる

### 正規分布の特徴

- 分布の中心は平均値  $\mu$  で、最も高い値(極大値)をとる
  - 平均値が変化すると、分布が左右に移動する
- $\mu \pm \sigma$  で変曲点(曲線の凹凸の変わり目)になる
  - 標準偏差が変化すると、分布の高さや広がりが変化する
- 平均値、中央値、最頻値は一致する
- 平均値と標準偏差から、分布の割合(曲線とx軸に囲まれる面積)が決まる
  - $\mu \pm \sigma$  の範囲 : 全体の約 68.24% を含む
  - $\mu \pm 2\sigma$  の範囲 : 全体の約 95.44% を含む
  - $\mu \pm 3\sigma$  の範囲 : 全体の約 97.73% を含む
  - それ以外の範囲 : 全体の約 0.27% を含む

## 標準正規分布

### 標準正規分布とは

- 正規分布で、平均値が 0、標準偏差が 1 になるように、正規分布の確率密度関数の変数  $x$  を次のように変換する(変数変換)

$$z = \frac{x - \mu}{\sigma}$$

- すると、標準偏差の式は次のように変形される

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

- 確率密度関数が上の式になるとき、『 $x$  は、標準正規分布にしたがう』という
  - 『 $z$  は、平均が 0 で標準偏差が 1 の正規分布にしたがう』ことになり、 $N(0,1^2)$  と表すことができる
  - つまり、変数だけから分布が決まる
  - さまざまな推定や検定に使われる

### 標準化 (基準化) (standardization)

- 先ほどの変数変換を、『標準化』または『基準化』という

$$z = \frac{x - \mu}{\sigma}$$

- (感覚としては) 分布を平均値の分だけ 0 まで移動し、分布の広がりを  $\frac{1}{\sigma}$  にする
  - 単位の異なるデータや平均値・分散が異なるデータを比較するときに使う (英語と数学のテストの成績の比較)
- ちなみに**偏差値**は、標準化を応用したもので、次のような式になる

$$z = 50 + 10 \times \frac{x - \mu}{\sigma}$$

### 標準正規分布の特徴

- 基本的な特徴は、正規分布と同じ
- 分布の曲線とx軸で囲まれた全体の面積は1になり
  - すなわち、面積が確率(割合)を表すことになる

## カイ二乗分布

### カイ二乗分布とは

- 変数  $z_1, z_2, \dots, z_r$  が標準正規分布にしたがい、互いに独立であるとする
- 次のようになるとき、『 $z$  は、自由度  $r$  のカイ二乗 ( $\chi^2$ ) 分布にしたがう』という

$$\begin{aligned} \chi^2 &= z_1^2 + z_2^2 + \dots + z_r^2 \\ &= \sum_{i=1}^r z_i^2 \\ &= \sum_{i=1}^r \left( \frac{x_i - \bar{x}}{\sigma} \right)^2 \end{aligned}$$

- 自由度とは分布の形状に影響を及ぼす値で、自由度の値が変わると分布の形状も変化する

### カイ二乗分布の特徴

- 平均は  $r$ 、分散は  $2r$  になる
- 適合度や独立性の検定など、分散の推定や検定に利用される
- 自由度が大きくなると、対称な分布に近づく

## t分布

### t分布とは

- 標準正規分布にしたがう確率変数  $x$  と、自由度  $r$  のカイ二乗分布  $X_r^2$  にしたがう確率変数  $y$  があり、互いに独立であるとする
- 次のようになるとき、『 $x$  は、自由度  $r$  のt分布にしたがう』という

$$t = \frac{x}{\sqrt{y/r}} = \frac{x}{\sqrt{\frac{X_r^2}{r}}}$$

- (感覚としては)正規分布にしたがう統計量を標準化したものの分布を表す
- この分布を発表した William Gosset のペンネームから「スチューデント(Student)のt分布」とも呼ぶ

### t分布の特徴

- 平均は0、分散は  $\frac{r}{r-2}$  になる
- 自由度が大きくなると、標準正規分布に近づく
  - 自由度が小さいときは、正規分布よりも裾の長い分布になる
- 母平均の推定や検定、平均値の差の検定や検定など、多くの統計的推定で利用される

## F分布

### F分布とは

- 確率変数  $x$  と確率変数  $y$  が、それぞれ自由度  $m$  と  $n$  のカイ二乗分布にしたがい、互いに独立であるとする
- 次のようになるとき、『 $x$  は、第一自由度が  $m$  で第二自由度が  $n$  のF分布にしたがう』という

$$F = \frac{x/m}{y/n} = \frac{X_m^2/m}{X_n^2/n}$$

- (感覚としては)分散の比率についての分布を表す

### F分布の特徴

- 2つの自由度を持つ
- 分散検定や分散分析(分散比の分布を調べる)に利用される
- 平均は  $\frac{n}{n-2}$ 、分散は  $\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}$  になる

# 母集団と標本

## 母集団と標本

### 全数調査と標本調査

- 全数調査、または悉皆(しっかい)調査
  - 全体について調べる(例:国勢調査)
- 標本調査
  - 全体から選び出した一部分について調べる(例:選挙予測)
  - 時間的・経済的に現実的な調査

### 母集団と標本

- 母集団(population)
  - 標本調査のもとになる集団
  - 母集団の大きさを、 $N$  であらわす
- 標本(sample)
  - 母集団から抽出(サンプリング)された一部分
  - 標本の大きさ(サンプルサイズ:標本に含まれる個数)を、 $n$  であらわす

## 標本の抽出

- 推測統計学(inductive statistics)
  - 標本にもとづいて母集団の特性値(母数:平均、分散などのパラメータ)を推定・予測する
  - 2つ以上の母集団の特性値を比較・検討する
- 標本が「母集団の精巧なミニチュア」になるように、偏りなく標本を抽出しないといけない

### 無作為抽出法(random sampling)

- 母集団から無作為に標本を抽出する
- 乱数表(無規則に数字を羅列した表)等をもとに、データ(個体)を標本として選ぶ
- どの標本も全く等しい確率で選ばれる(主観や作為などが入る余地がない)
- 母集団が大きい場合は実施が困難(例えば数万人の名簿を作るのは容易ではない)

### 系統抽出法(systematic sampling)

- 最初の標本のみ乱数表などで選ぶ
- あとの標本は一定間隔(抽出間隔;sampling interval)で抽出する
- 実際には抽出間隔が扱いやすい数でない場合が多い
- 無作為抽出法と同じで、母集団が大きい場合は実施が困難

### 多段抽出法(multi-stage sampling)

- 何段階かのステップを経て標本を抽出する(例:2段階抽出法)
- 何段階かのステップで標本にする対象を絞り込み、最後に無作為抽出法で標本を選ぶ
  - 4段階の例:全国 都道府県 市町村 病院・診療所 看護師(あとは無作為抽出)
- 調査の手間ひまが少なく実践的だが、標本の偏り(バイアス)に注意しないといけない

### 層別抽出法(stratified sampling)

- 母集団を同じ特徴(年代、職業、都道府県など)で層に分ける(層別化)

- 層の構成比に応じて適切な数だけ、層ごとに無作為抽出法で標本を選ぶ
- 層の違いによって偏りがある状況が事前に予測される場合に適している (地域と政党の支持率、年代と好きな芸能人)
  - 母集団の特徴を少なくとも1つは事前に知っている必要がある