

母集団の統計的推定

心理データ解析演習 2011/07/06

M1 水垣 さなえ



発表の流れ

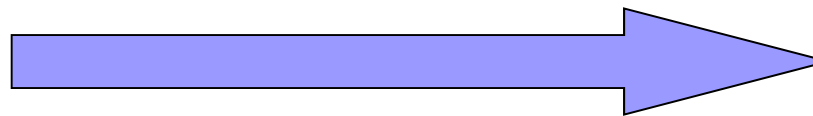
- はじめに
- 推測統計の基礎事項
- 母集団平均の区間推定 (t 分布等)
- 母集団分散の区間推定 (χ^2 分布)
- 母集団比率の区間推定 (2項分布)
- 区間推定の注意

はじめに 統計的推定とは

例題(管,2009より)

: Aさんは久しぶりに家に帰り、稲刈りをしました。収穫された稲穂を1本拾い粒数をかぞえると90粒でした。別の1本についても調べると95粒でした。さらに別の稲穂を抜き出し全部で20本の粒数を数えました。稲穂1本あたりの平均粒数は93粒であることがわかりました。このことからAさんは次のことを考えました。

1. 実家の田んぼ全体の稲穂一本あたりの平均粒数を、たった一度の調査である20本の平均値から93粒と判断してよいだろうか。
2. 1つの値で言い切るのが難しければ、90粒から95粒の間にあるといった大雑把なことはいえないだろうか



この2つの問いに答える方法が

統計的推定

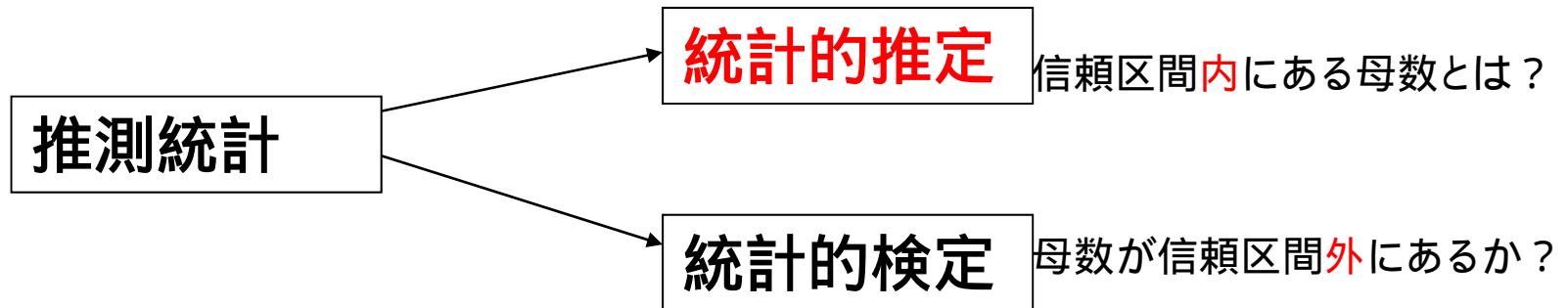




基礎事項

推測統計

調べたい集団の**一部**のデータから集団**全体**の特色や傾向を明らかにする方法。



統計的推定と統計的検定は”表裏一体”の関係

推計したい集団のこと

= 母集団

集団の一部を対象とする調査

= 標本調査

標本調査における対象

= 標本orサンプル

標本調査における対象の個数

= サンプルサイズ

推測統計の用語

推測統計では標本調査の基本統計量を用いて、母集団を推測する

母集団

- サイズ N
- 母平均値 μ (ミュー)
- 母標準偏差 (シグマ)
- 母分散 σ^2
- 母比率 p

標本調査

- 標本サイズ n
- 標本平均 \bar{x}
- 標本標準偏差 s
- 標本分散 s^2
- 標本比率 \bar{p}

推測統計の主な考え方(確率)

確率事象

ある確率変数(X)が出るというような事柄に目をつけたとき、その事柄が現れるのは、1回の施行では不確定であるが、その事柄の起こる割合が多数回の実験を試みるとき安定性をもつ、そういう事柄だけについて考えるとき、これを**確率事象**とよぶ。

確率事象は、標本調査の基本統計量間の関係から主に次のようなことが分かっている。

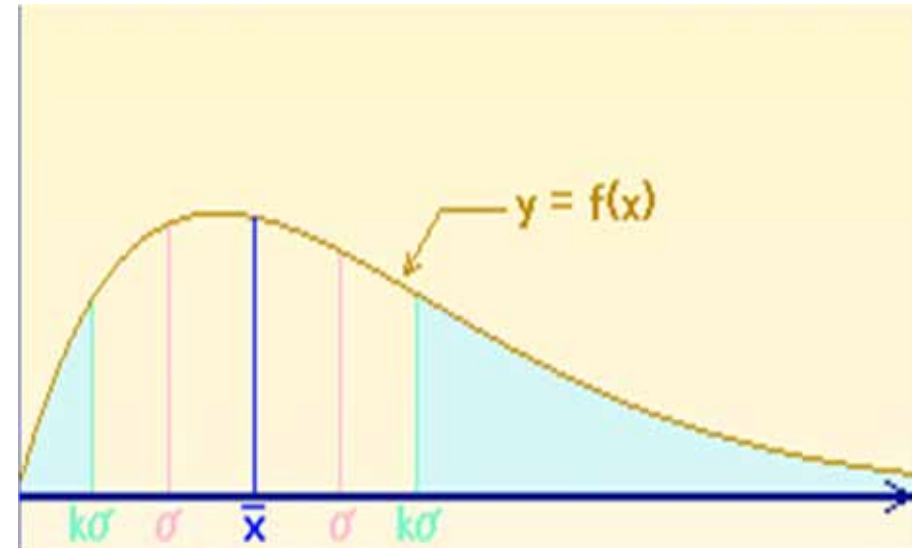
1. **チェビシェフの不等式**
2. **正規分布**
3. **大数の法則**
4. **中心極限定理**

統計的推定の主な考え方(確率)

チェビシェフの不等式

確率変数 X の平均 $E[X] = \mu$ 、分散 $V[X] = \sigma^2$ が共に有限ならば任意の $k (> 0)$ に対して

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$



<http://www.kwansei.ac.jp/hs/z90010/sugakuc/toukei/cebysev/cebysev.htm>より

つまり母平均 $\mu (= \bar{x})$ 以外の確率変数 X が出る割合(母平均 μ からのズレの割合)は、母標準偏差 σ の値で分かるというもの。

式の意味は、母平均 μ と他の確率変数の差が母標準偏差 σ の k 倍以上になる確率は $1/k^2$ 以下であるということ。(上の図の青い面積部分の確率)

統計的推定の主な考え方(確率)

正規分布

理論的に考案された分布

理論的な式によって正規分布に従う変数 X が任意の a から b までの値をとる確率 P

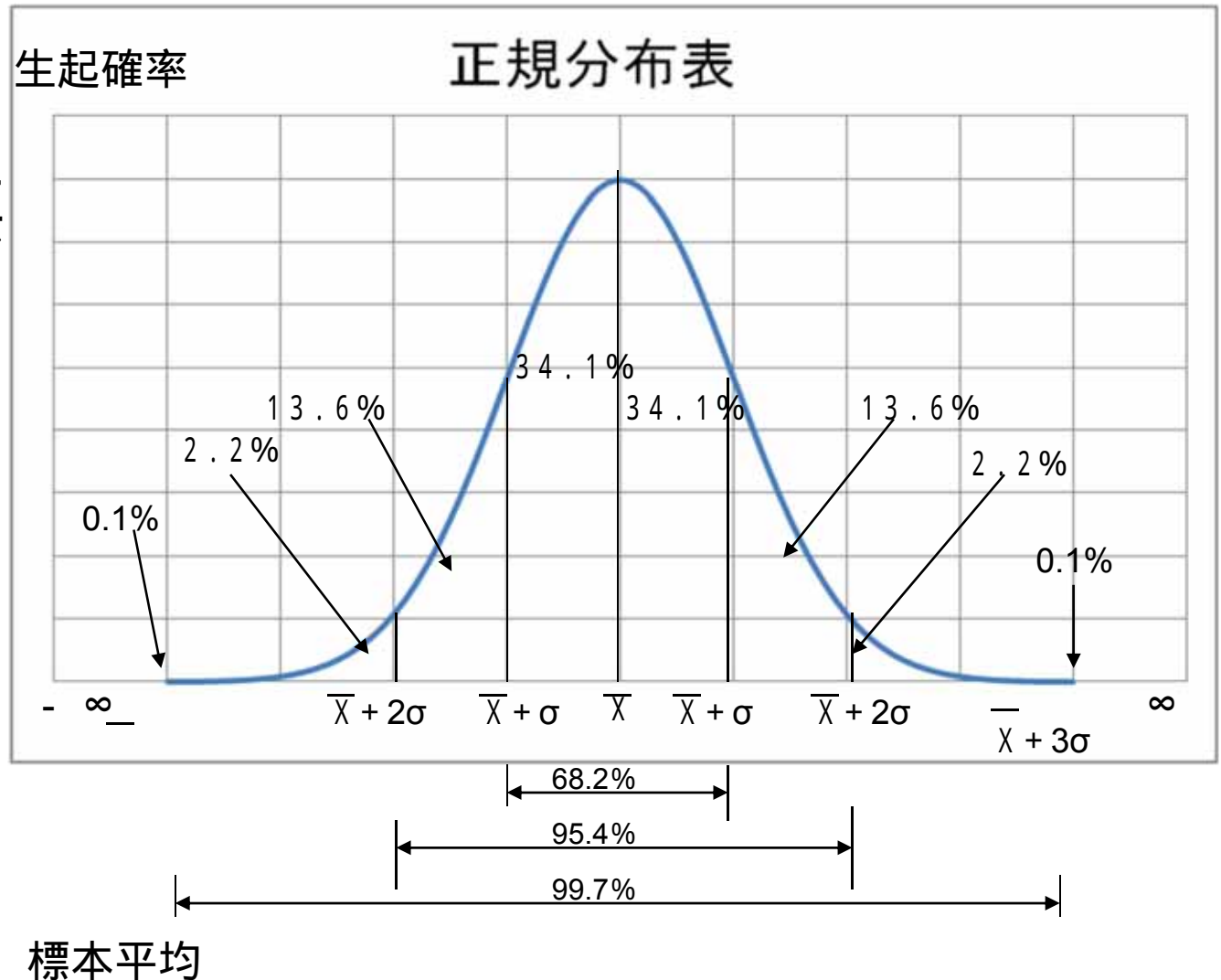
($a < X < b$)が分かる

左右対称(歪度 = 0)

中央に山が一つ

両すそがなだらかに広がった(尖度 = 3)形

ベルカーブ(釣鐘状)



統計的推定の主な考え方(確率)

大数の法則

どんなサンプルと母集団にも当てはまる法則

確率変数 X_1, X_2, \dots, X_n が独立であって、その期待値 $E(X_1) = E(X_2) = \dots = E(X_n) = \mu$ かつ、一定の定数 σ^2 に対して X_i の分散 $V(X_i) = \sigma^2 (i=1, 2, \dots, n)$ であるならば、任意に与えられた正数 $\epsilon > 0$ に対して、

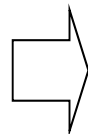
$$\lim_{n \rightarrow \infty} P \left(\left| \frac{X_1 + X_2 + \dots + X_n}{n} - \mu \right| < \epsilon \right) = 1$$

となる。

これは、 n を十分に大きくとれば確率変数 \bar{X} は定数 μ に限りなく近づくという意味



サイコロを一回
振って1の目が
出る確率 =
=理論的確率



出る目は1
か1以外で
しかなく、1
が1/6回で
るといこと
はありえな
い



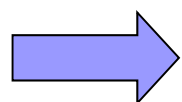
でも多数回
ふると1が出
る割合は頻
度の比1/6と
して確認でき
る。

統計的推定の主な考え方(確率)

中心極限定理

X が平均 μ 、標準偏差 σ のある分布に従うとき、大きさ n の無作為標本にもとづく標本平均 \bar{X} は、 n が無限に大きくなるとき、平均 μ 、標準偏差 σ / \sqrt{n} の正規分布に近づく

中心極限定理は**母集団の分布に関係なく**、




その標本平均は、正規分布に従う

<http://www.kwansei.ac.jp/hs/z90010/sugakuc/toukei/tyuusin2/chuusin.htm>



推定法



稲穂の例題(管,2009)にもどると

Aさんは久しぶりに家に帰り、稲刈りをしました。収穫された稲穂を1本拾い粒数をかぞえると90粒でした。別の1本についても調べると95粒でした。さらに別の稲穂を抜き出し全部で20本の粒数を数えました。稲穂1本あたりの平均粒数は93粒であることがわかりました。このことからAさんは次のことを考えました。

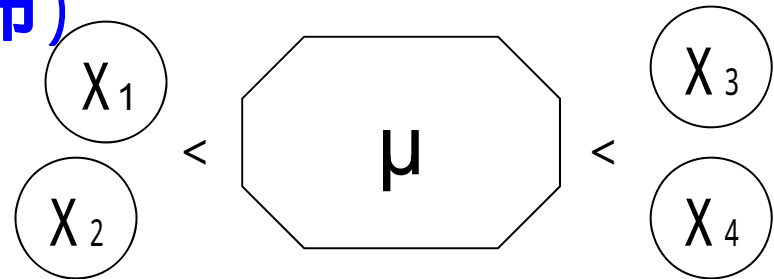
1. 実家の田んぼ全体の稲穂一本あたりの平均粒数を、たった一度の調査である20本の平均値から93粒と判断してよいだろうか。
2. 1つの値で言い切るのが難しければ、90粒から95粒の間にあるといった大雑把なことはいえないだろうか

標本分布と母数の関係

標本分布(標本平均の分布)

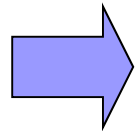
- 無作為の標本抽出を多数回それぞれ独立に繰り返したとき、標本抽出に伴う誤差(偶然)によって標本の平均値 \bar{X} は、母平均 μ と完全に一致することはほとんどなく、母平均 μ を中心にある程度のバラツキ(分散)をもって分布する。

= 標本分布(標本平均の分布)



稲穂の粒数の例でいえば...

もし20本の標本調査を何度も繰り返した時、最初の稲穂の粒数の平均は93粒だったが、別の調査では、95粒だったり、90粒だったりするかもしれない



標本分布と母数の関係

標本分布の分散

母集団が正規分布に従う時、**標本分布**も正規分布に従うことが知られており、その分散は

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$$

標本サイズ n によって標本分布の分散の大きさが変わる
= n が大きくなれば、よりバラツキの少ない、つまり母平均 μ の近くに標本の平均値(= 標本分布における確率変数)があつまると考えられる

標準誤差

標本から母集団について推定を行う時、標本分布の分散の大きさに応じて誤差が伴う。このため、標本分布の標準偏差を**標本誤差**(Standard Error: SE)と呼ぶ



母数の推定

標本分布と母集団の関係性から母平均 μ と母分散 σ^2 の値を推定していく

- 1) 不偏推定値による母数の推定
- 2) 不偏推定値をもとにした母数の区間推定

不偏推定値による母数の推定

不偏推定値

ある統計量の期待値がそれに対応する母数(母集団の統計量)と一致する場合、その統計量を母数の**不偏推定値**という。

1) 母平均の推定

大数の法則や**中心極限定理**などの考え方から、

標本平均 \bar{x} は母平均 μ の**不偏推定値**ということができる。

稲穂の例題 問1への答え

「実家の田んぼ全体の稲穂一本あたりの平均粒数を、たった一度の調査である20本の平均値から93粒と判断してよいだろうか。

YES

不偏推定値による母数の推定

2) 母分散の推定

母分散 σ^2 も標本分散 s^2 をそのまま不偏推定値にできるか？

母分散 σ^2 は、



標本の分散 s^2
(\bar{X} に対する確率変数 X のバラツキ具合)
+
標本分布の分散 $\sigma_{\bar{X}}^2$
(μ に対する \bar{X} のバラツキ具合)

だから、母分散 σ^2 s^2 で
母分散の不偏推定値は $\hat{\sigma}^2$

$$= \frac{\sum_i^n (X_i - \bar{X})^2}{n} + \frac{\sigma^2}{n}$$

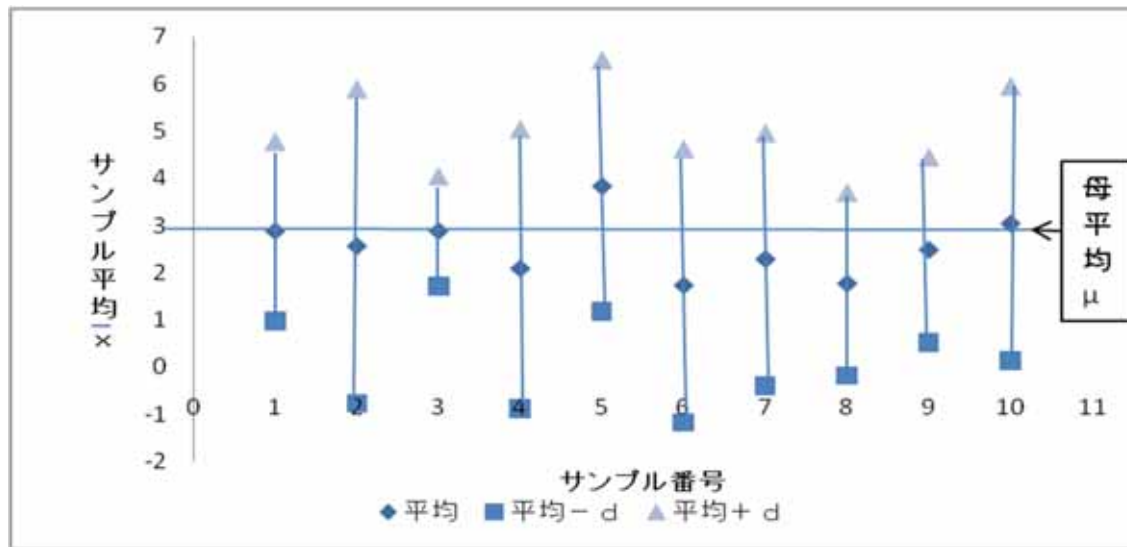
$$= \frac{\sum_i^n (X_i - \bar{X})^2}{n - 1}$$

$$= \hat{\sigma}^2$$

不偏推定値をもとにした母数の区間推定

区間推定法

標準誤差(標本抽出による誤差)を考慮して、標本平均 \bar{X} を含むひとつの区間を設定し、その区間の中に母平均 μ があるようにする方法



林、1991を参考に作成

稲穂の例題 問2への答え
1つの値で言い切るのが難しければ、90粒から95粒の間にあるといった大雑把なことはいえないだろうか

YES

不偏推定値をもとにした母数の区間推定

信頼区間

“母平均 μ が区間～から～までの間に存在する確率は95%である”と推定の信頼性が確率的に表現された区間

信頼限界

推定された母数の上限および下限

たとえば母平均 μ の推定の信頼区間の上限(\bar{X}_u)および下限(\bar{X}_L)は、

上限 $\bar{X}_u = \bar{X} + \text{誤差}$ 下限 $\bar{X}_L = \bar{X} - \text{誤差}$
誤差の値の作り方は後述

信頼係数(信頼度)

推定の信頼性を表す確率。

例えば母平均推定の場合、母平均 μ が信頼区間に含まれる確率を意味する。この場合信頼係数は $1 - \alpha$

誤差と信頼度は裏腹の関係にある。 誤差が大(信頼区間広い) 信頼度が低い
誤差が小(信頼区間狭い) 信頼度が高い

母数の区間推定

平均の推定

■ 母集団平均の推定 (1)

母集団が正規分布 $N(\bar{X}, \sigma^2)$ で σ^2 が既知の時

標準正規分布を用いて推定 (実際の調査ではほとんどないが考え方として)

$$P_r \left\{ |\bar{X} - \mu| < k_\alpha \frac{\sigma}{\sqrt{n}} \right\} = P_r \left\{ \frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} < k_\alpha \sigma \right\}$$

= ((t) の $-k_\alpha$ から k_α までの面積)

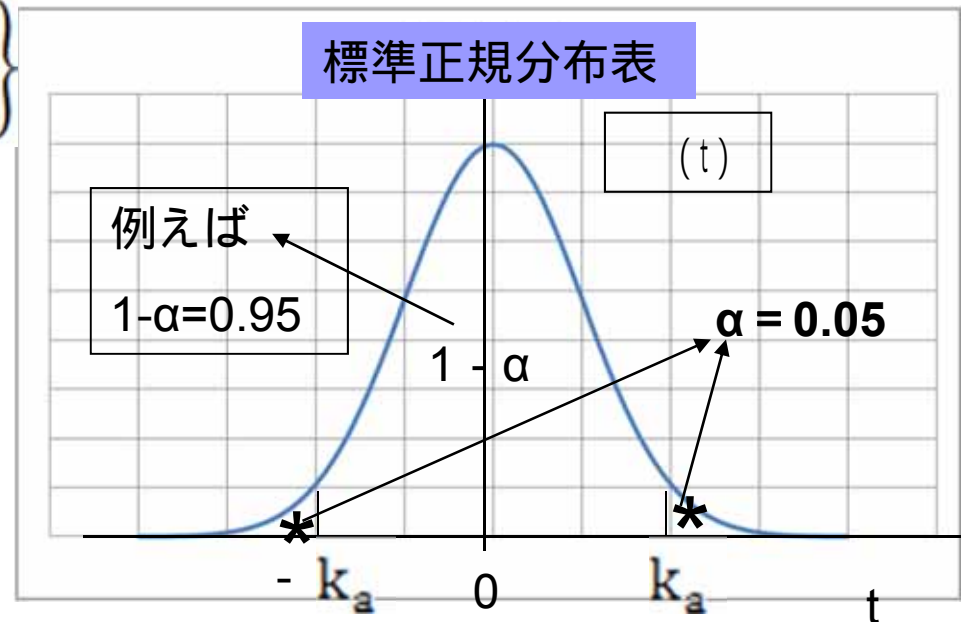
= 1 -

たとえば信頼度を95%とするなら、

$1 - \alpha = 0.95$ で $\alpha = 0.05$ となる z 値を

標準正規分布表から探して計算

$$\bar{x} - \frac{1.96\sigma}{\sqrt{n}} < \mu < \bar{x} + \frac{1.96\sigma}{\sqrt{n}}$$



母数の区間推定

平均の推定

■ 母集団平均の推定 (2)

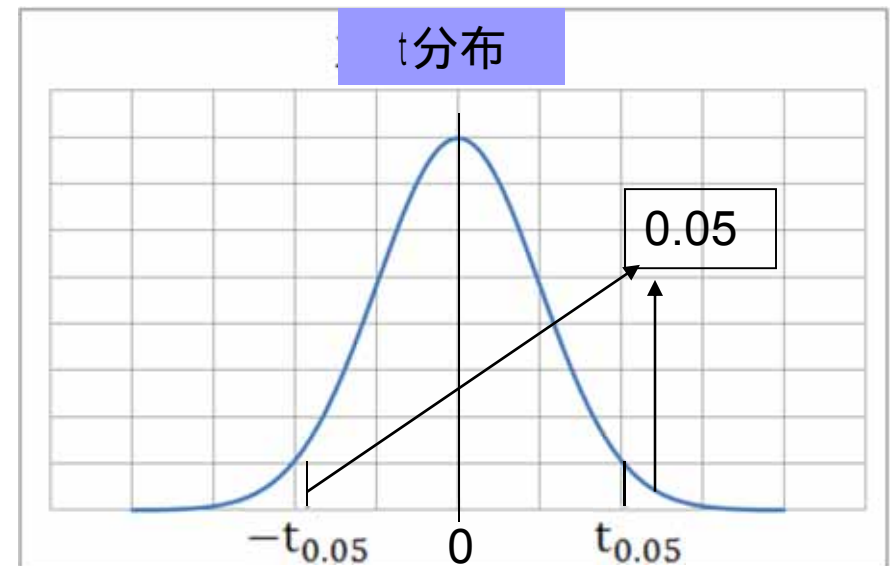
母集団が正規分布 $N(\bar{X}, \sigma^2)$ で σ^2 が未知の時

t分布 を用いて推定 (自由度は $n-1$) 信頼係数 95% の時

$$T = \frac{\sqrt{n-1}(\bar{x} - \mu)}{s} \quad \text{とすると}$$
$$P_r\{-t_{0.05} < T < t_{0.05}\} = 0.95$$

となる t の値は $t_{0.05} = 2.776$

$$= \bar{X} - t_{\alpha} \frac{s}{\sqrt{n-1}} < \mu < \bar{X} + t_{\alpha} \frac{s}{\sqrt{n-1}}$$
$$= \bar{X} - 2.776 \frac{s}{\sqrt{n-1}} < \mu < \bar{X} + 2.776 \frac{s}{\sqrt{n-1}}$$



母数の区間推定

平均の推定

母集団平均の推定(3)

■ 母集団の分布が不明の時

通常の調査では標本平均 \bar{X} について、 n が大きい時(実用上は大体30以上)

$$E(\bar{X}) = \mu, \quad \sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$$

なる正規分布に近似的に従うため、**標準正規分布**で推定
また n が十分大きい時(大体100以上)は

標本分散 s^2 を母平均 σ^2 の代用としてもよい(**中心極限定理**)。

信頼係数は100%としてもよいが、現実的に95%として

$$\bar{x} - \frac{1.96s}{\sqrt{n}} < \mu < \bar{x} + \frac{1.96s}{\sqrt{n}}$$

母数の区間推定 分散の推定

母集団分散の推定

- 母集団が正規分布 $N(\bar{X}, \sigma^2)$ で σ^2 が未知の時
 χ^2 分布を用いて推定 (自由度 $n-1$)

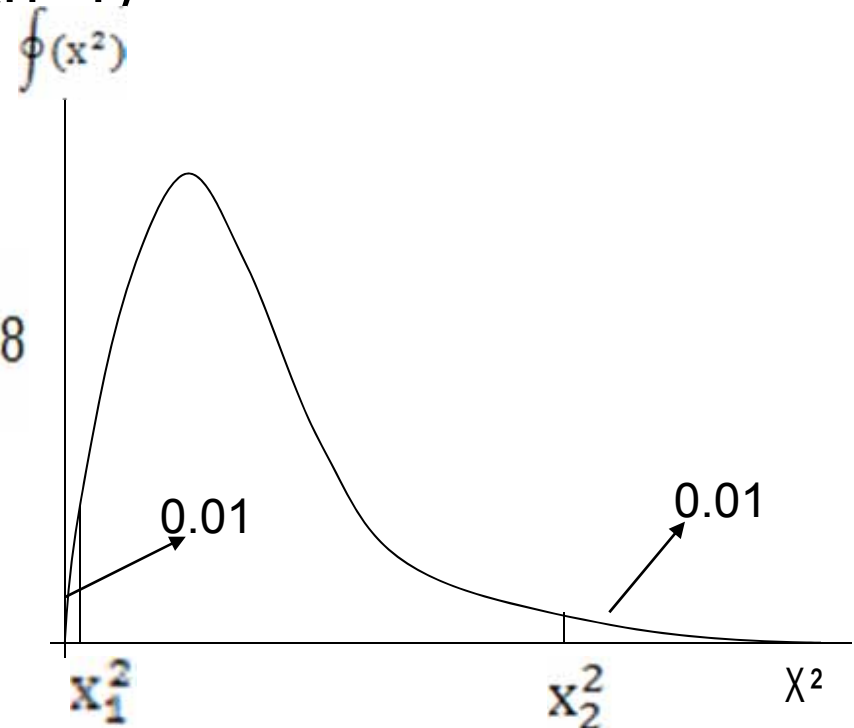
$$\chi^2 = \frac{ns^2}{\sigma^2} \quad \text{とすると}$$

たとえば信頼係数98%の時

$$P_r\left\{x_1^2 < \chi^2 < x_2^2\right\} = P_r\left\{\frac{ns^2}{x_2^2} < \sigma^2 < \frac{ns^2}{x_1^2}\right\} = 0.98$$

を満足する x_1^2 と x_2^2 を χ^2 分布より求めて

$$\frac{ns^2}{13.277} < \sigma^2 < \frac{ns^2}{0.297}$$



母数の区間推定 比率の推定

母集団比率の推定

■ 母集団が二項分布の時

$E(p)$ 標本比率の期待値 = $E(\bar{X})$ 標本平均の期待値 = P 母比率

このとき母集団の分散 = $P(1-P)$

標本の大きさが十分大きい時(大体100以上)は標本比率は p は近似的に平均 P , 分散 $\frac{P(1-P)}{n}$ の正規分布に従うから

$$P_r \left\{ |p - P| < k_\alpha \sqrt{\frac{P(1-P)}{n}} \right\} = 1 - \alpha \quad \text{より}$$

信頼係数 $1-\alpha$ の信頼区間は

$$p - k_\alpha \sqrt{\frac{P(1-P)}{n}} < P < p + k_\alpha \sqrt{\frac{P(1-P)}{n}}$$



区間推定の注意点



区間推定の注意点

確率モデルを適用するためには

1. 標本調査が無作為(ランダムサンプリング)
2. どのような母集団を推定するか明確に
3. 母集団分布が正規分布に従うことを仮定
4. 推定によって得られるのは数学的な結果

参考文献

- 管 民郎(2009)『らくらく図解 統計分析教室』, (株)オーム社
- 丹羽時彦 関西学院高等部ホームページより
「放課後の数学入門」 統計(数学C)
<http://www.kwansei.ac.jp/hs/z90010/sugakuc/toukei/toukeihy.htm>
- 林 知己夫編 林 文、佐藤良一郎、青山博次郎、林知己夫著(1991)
『統計学の基本』 朝倉書店
- 吉田寿夫 (1998)
『本当にわかりやすいすごく大切なことが書いてあるごく初歩の統計の本』 (株)北大路書房
- 森 敏昭、吉田寿夫編著(1990) 『心理学のためのデータ解析テクニカルブック』
(株)北大路書房