

スライド2 データの整理（1）

1. データの中心の代表値

1.1 データを元にした平均，中央値，最頻値の計算法

1.2 平均，中央値，最頻値の特徴

- ・ 平均
 - ・ 極端な値（はずれ値，異常値）の影響大
 - ・ 計算が容易
- ・ 中央値
 - ・ はずれ値（異常値）の影響を受けない
- ・ 最頻値
 - ・ はずれ値（異常値）の影響を受けない
 - ・ 最頻データ以外を完全に無視している

2. 散らばりの代表値

2.1 分散，標準偏差，範囲

- ・ 標本分散（標本が調査対象の一部の場合）

$$\frac{(\text{データ}-\text{平均})^2 \text{の合計}}{\text{標本数}-1}$$

- ・ 分散（標本が調査対象と一致する場合）

$$\frac{(\text{データ}-\text{平均})^2 \text{の合計}}{\text{標本数}}$$

- ・ (標本)標準偏差 $\sqrt{(\text{標本})\text{分散}}$

- ・ 範囲 = 最大値 - 最小値

2.2. チェビシェフの不等式

- ・ チェビシェフの不等式（どんなデータにも当てはまる）

$$|\text{データ値}-\text{標本平均}| \geq k \times \text{標本標準偏差} \text{を満たすデータの割合は全標}$$

本数の $1/k^2$ 以下である．

2.3. その他散らばりの代表値

- ・ 四分位点（しぶんいてん）

第1四分位点

- ・ 順序どおりに並べ替え，データを4分割する．その4分割点のうち最も小さい値．

第2四分位点 中央値と同じ

- ・ 現在のデータに比重をおく
- ・ 単純平均ではなく加重平均をとる

6. 変動係数

- ・ 散らばりを比較する

平均が同じくらい

- ・ 分散や範囲，四分位範囲を比較すればよい

平均が異なるデータの比較

- ・ 平均で標準化した標準偏差 = 変動係数 $\frac{\text{標本標準偏差}}{\text{標本平均}}$

- ・ 四分位分散係数

四分位範囲の半分を中央値で割ったもの

- ・ 十分位分散係数 $\frac{\text{第9十分位数} - \text{第1十分位数}}{2 \times \text{中央値}}$

スライド3 データの整理(2)

1. ローレンツ曲線

- ・ (正の) 経済変量の集中度を表すのに使用される

- ・ データ値による描き方

データ値を小さい順に並べる。

各データについて，(順位/標本数)と(それ以前のデータに関する経済変量の合計/全データに関する合計)の組を求めておく。

この組み合わせを(相対順位，相対変量)と呼ぶ

その組み合わせに対応した点をプロットし，折れ線グラフを書く。

- ・ 度数分布表による描き方

各階級について，(累積相対度数)と(その階級までの<階級値×度数>の合計/全データに関する<階級値×度数>の合計)の組を求めておく。

その組み合わせに対応した点をプロットし，折れ線グラフを書く。

2. ジニ係数

- ・ 相対順位

あるデータ値の順位(小さい順)のデータ数に対する割合

度数分布表の場合ある階級の累積相対度数

- ・ 相対変量

それ以下のデータ値の合計を全てのデータの合計で割った値

度数分布表の場合ある階級の(階級値×相対度数)/全データの合計

全データ値の合計は(階級値×相対度数)の合計

- ・ 定義

$$G = \left\{ \left(\text{各階級とその一つ前の相対順位} \right) \times \text{相対変量} \right\} - 1$$

図形的意味

右図のように各階級毎にその相対順位と相対累積所得に組み合わせに点を打っていき、それつなぎ合わせるとローレンツ曲線ができる。

ジニ計数値は45度線とローレンツ曲線との間の面積の2倍

極端なローレンツ曲線とそのときのジニ係数

すべてが同一値、または、同一階級の場合（平等ケース）

・ ジニ係数は0

1 データだけが極端に大きい場合

・ ジニ係数は1

ジニ係数の特徴

・ 端のデータが真ん中に移ると、ジニ係数は小さくなる（だから不平等の指標）

・ データの単位を変えてもジニ係数は変わらない

3.2 変数データの整理

3.1 散布図

・ 2変数の関係を視覚的に捉える = 散布図

データ値の組み合わせを座標と考えて、データ毎に点をプロットする。

3.2 標本相関係数

・ 2変数の関係の深さを数値化する。

・ 共分散

$$\text{標本共分散} = \frac{\left\{ \begin{array}{l} (x\text{変数の値} - x\text{の標本平均}) \times (y\text{変数の値} - y\text{の標本平均}) \\ \text{のすべての標本に関する合計} \end{array} \right\}}{n(\text{または}n-1)}$$

・ 分散を計算するのにnで割ったかn-1で割ったかによってきめる

共分散の性質

・ xの値がその平均を上回るときにはyもその平均を上回る様な場合、共分散の式の分子は正になる（右上がりの傾向）

・ xの値がその平均を上回るときにはyはその平均を下回る様な場合、共分散の式の分子は負になる（右下がりの傾向）

・ xの値とyの値に関係がない場合は、0になる

・ 標本相関係数

$$r_{xy} = \frac{\text{標本共分散}}{x\text{の標本標準偏差} \times y\text{の標本標準偏差}}$$

・ 標本相関係数の性質

符号は相関の方向を表す

- ・ x が平均より大きかったら y も平均より大きい傾向
符号は正，正の相関を持つと呼ぶ

- ・ x が平均より大きかったら y は平均より小さい傾向
符号は負，正の相関を持つと呼ぶ

絶対値は相関の強さ，関係の深さを表す

- ・ 大きいほど関係は深い
- ・ 0 ~ 1 の間

標本相関係数は - 1 以上 + 1 以下

相関係数の限界

相関係数は線形関係が完全に成立している場合を基準にしている < 相関係数を計算する前に散布図を書こう >

- ・ 非線形関係，例えば，2 次式の関係などは表せない。

見せかけの相関

- ・ 二つの変数とも，第 3 の変数と相関を持っていて，その結果，両変数に関係があるように見える場合がある。
- ・ 例えば，時系列データの場合に，二つの変数が時間とともに増加，または，減少している場合

3.3 標本偏相関係数

- ・ 第 3 の変数と両方とも相関を持っていて，そのために 2 つの変数の相関が強いように見える場合の対策

第 3 の変数 (z とする) の影響を取り除いたうえで，二つの変数の相関を調べる。

標本偏相関係数
$$r_{xy|z} = \frac{r_{xy} - r_{xz}r_{zy}}{\sqrt{(1-r_{xz}^2)(1-r_{zy}^2)}}$$

スライド 4 確率

0. 確率

- ・ 事象 出来事
- ・ 事象 (出来事) がおきる可能性の指標
確実に起きるが 1 ，起きないが 0
起きる可能性が確実に起きる場合の何%か？
- ・ 様々な事象を表すために A ， B など大文字を使用
A が起きる確率を P (A) で表す
- ・ 積事象：どちらも起きるといふ出来事 $A \cap B$
- ・ 排反事象

AかBのどちらかしか起きないという関係

AとBと一緒に起きるといことはない

和事象：どれかが起きるとい出来事 $A \cup B$

これはAとBのどちらかしか起きないという意味ではない

もし、AとBが排反事象なら、 $P(A \cup B) = P(A) + P(B)$

余事象

「ある出来事が起きない」とい出来事

Aの余事象 A^c で表す

$$P(A^c) = 1 - P(A)$$

$P(A \cup B) = P(A) + P(B) - P(A \cap B)$

1. 条件付き確率

1.1 条件付き確率の概念 (1)

2段構えの確率 = 条件付き確率

つまりある情報が解っている状況に限定し、そのもとで事象が起きる確率

1.2 条件付き確率と同時確率

$$P(X \cap Y) = P(X | Y) \times P(Y) = P(Y | X) \times P(X)$$

1.3 条件付き確率と普通の確率

$$P(B_1)$$

$$= P(B_1 \cap A_1) + P(B_1 \cap A_2) + P(B_1 \cap A_3)$$

$$= P(B_1 | A_1)P(A_1) + P(B_1 | A_2)P(A_2) + P(B_1 | A_3)P(A_3)$$

1.4 同時確率から条件確率へ

$$P(X | Y) = P(X \cap Y) / P(Y)$$

$$P(Y | X) = P(X \cap Y) / P(X)$$

2. 独立

2.1 独立性の定義

$$P(X | Y) = P(X), \text{ 同じことだが, } P(X \cap Y) = P(X | Y) \times P(Y) = P(X) \times P(Y)$$

3. ベイズの定理

$$P(A_i | B_k)$$

$$= \frac{P(A_i \cap B_k)}{P(B_k)} = \frac{P(B_k | A_i)P(A_i)}{P(B_k | A_1)P(A_1) + P(B_k | A_2)P(A_2) + P(B_k | A_3)P(A_3)}$$

スライド5 確率分布

1. 確率変数

事象を数値化したもの

(事象 → 数値) の関数

2. 離散確率変数と連続確率変数

3. 確率関数

3.1 概念

- 飛び飛びの値それぞれになる確率を示す

確率関数

- $X = v_i \quad i=1, \dots, k$, つまり, 確率変数 X は k 個の飛び飛びの値をとるとする.

$$p(t) = p_X(t) = \begin{cases} t = v_i \text{ の場合} & P(X = v_i) \\ \text{それ以外の場合} & 0 \end{cases}$$

4. 確率分布関数

4.1 概念

- 離散確率変数でも連続確率変数でも定義可
- 確率変数 X の分布関数 $F(t) = F_X(t) = P(X \leq t)$
- 離散確率変数の場合

t 以下の値をとる確率の合計

$$F(t) = \sum_{i=1}^k p(v_i) \times I_{v_i \leq t}(t) = \sum_{\substack{i=1 \\ v_i \leq t}}^k p(v_i)$$

4.3 確率分布関数と区間確率

$$P(a < X \leq b) = P(X \leq b) - P(X \leq a) = F(b) - F(a)$$

5. 確率密度関数

連続確率変数の場合 $P(X = t) = 0$ になるので確率関数は定義できない.

$$f(t) = f_X(t) = \lim_{\substack{a \rightarrow t \\ b \rightarrow t}} \frac{P(a < X \leq b)}{b - a} = \lim_{\substack{a \rightarrow t \\ b \rightarrow t}} \frac{F(b) - F(a)}{b - a} = F'(t)$$

5.4 確率と確率密度関数

$$P(a < X \leq b) = P(a \leq X \leq b) = \int_a^b f(t) dt$$

$$P(X \leq b) = \int_{-\infty}^b f(t) dt = F(b)$$

離散の場合, 確率密度関数は存在しない.

6.5 確率関数と密度関数の基本性質

- 確率関数の場合

$$\text{確率関数の合計は } 1 \quad \sum_{i=1}^k f(v_i) = 1$$

- 確率密度関数の場合

密度関数の面積合計も 1

$$\int_{-\infty}^{\infty} f(t)dt = 1$$

6 . 分布の代表値

6 . 1 分布の平均 (母平均, 期待値)

- ・ 分布の重心
- ・ 計算法

離散確率変数の場合

$$E[X] = \mu = \mu_X = \sum_{i=1}^k v_i P(X = v_i) = \sum_{i=1}^k v_i p(v_i)$$

連続確率変数の場合

$$E[X] = \mu = \mu_X = \int_{-\infty}^{\infty} t f(t)dt$$

6 . 2 分布の分散 (母分散)

- ・ 確率分布の散らばりの指標
- ・ 計算法

離散確率変数の場合

$$V[X] = \sigma_X^2 = \sum_{i=1}^k (v_i - \mu_X)^2 p(v_i)$$

連続確率変数の場合

$$V[X] = \sigma_X^2 = \int_{-\infty}^{\infty} (t - \mu_X)^2 f(t)dt$$

6 . 3 確率変数から新たな確率変数を作る

- ・ 確率変数 X の関数もまた確率変数
- ・ この確率変数 $g(X)$ の分布関数は, $F_{g(X)}(t) = F_X(g^{-1}(t))$
- ・ 期待値計算

- ・ (離散) $E[g(X)] = \sum_{i=1}^k g(v_i)p(v_i)$

- ・ (連続) $E[g(X)] = \int_{-\infty}^{\infty} g(t)f(t)dt$

6 . 3 期待値, 分散の演算

期待値の性質 $E[X - \mu_X] = 0$

- ・ 期待値の演算

X, Y は確率変数, a, b は確率変動しないとする

$$E[aX + bY] = aE[X] + bE[Y]$$

- 分散の演算

$$V(aX + b) = a^2V(X)$$

X と Y が独立の場合

$$V(aX + bY) = a^2V(X) + b^2V(Y)$$

6.4 分布のパーセント点

- 確率変数 X の分布の %点

$$F(t) = P(X \leq t) = \alpha/100 \text{ となる } t \text{ の値}$$

- 分布の中央値 (メジアン)

$$F(t) = P(X \leq t) = 0.50 \text{ となる } t \text{ の値}$$

7. 正規確率変数と正規分布

7.1 独立な変数の和の分布

- 独立な確率変数の和の分布を考える (中心極限定理)

X_1, X_2, \dots, X_n を独立で期待値 $E[X_i] = 0$, 分散 $V(X_i) = 1$ の確率変数の列とする

$$S_n = X_1 + X_2 + \dots + X_n$$

S_n/\sqrt{n} の極限分布

- 標準正規分布とよぶ (N(0,1)と書く)

$$\text{密度関数 } f(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$$

- このような分布を持つ確率変数を Z とする.

一般の正規分布

- 平均 μ , 分散 σ^2 (標準偏差) の正規分布 ($N(\mu, \sigma^2)$ とかく) は確率変数 $\sigma Z + \mu$ の分布

$$\text{密度関数 } f(t) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}}$$

7.2 正規分布表

- 逆に平均 μ , 標準偏差 σ の正規分布は標準正規分布の確率変数 Z を使って $\sigma Z + \mu$

で表せるから, $\frac{X - \mu}{\sigma} = Z$ は N(0,1) の分布となる.

$$\text{従って, } P(X \leq s) = P\left(\frac{X - \mu}{\sigma} \leq \frac{s - \mu}{\sigma}\right) = P\left(Z \leq \frac{s - \mu}{\sigma}\right)$$

- 教科書 p. 279 の標準正規分布表を使えば計算可
- 正規分布は平均の値を軸にして左右対称な密度関数を持つ

- ・ 故に，平均の値を軸にして左右対称な確率分布関数を持つ
 - ・ よって， $P(Z \leq a) = P(-Z \leq a) = P(Z \geq -a)$
 - ・ または， $P(Z \leq -a) = P(-Z \geq a) = P(Z \geq a)$

7.3 偏差値の意味

偏差値 $S = \frac{X - \mu}{\sigma} \times 10 + 50$

偏差値は仮定の下で平均 50，標準偏差 10 の正規分布

従って， $P(S \leq t) = P\left(\frac{S-50}{10} \leq \frac{t-50}{10}\right) = P\left(Z \leq \frac{t-50}{10}\right)$

スライド 6 標本分布と推定

1. ランダムサンプリング

- ・ 無作為抽出法
 - ・ 調査対象のすべてを調査せず，その一部を無作為に選び，そのみを調査する
 - ・ 無作為に選ぶ場合，多くは乱数表を使用し，乱数表の示す番号に該当する対象を標本とする（単純無作為抽出）。

1.1 無作為抽出とは何か？

- ・ 無作為かどうかは，標本の選び方と，調査する属性が独立であるかどうかによって決まる。
- ・ 標本同士は独立でなければならない
- ・ 無作為抽出ではない場合を有意抽出と呼ぶ

2. 標本平均の分布

2.1 標本平均の期待値

- ・ 標本平均の期待値

$$E[\bar{X}] = E\left[\frac{X_1 + \dots + X_n}{n}\right] = \mu$$

標本平均の期待値は調査対象の分布の平均と一致（調査対象の母平均）

2.2 標本平均の（分布の）分散

- ・ 標本平均の分散

$$V[\bar{X}] = V\left[\frac{X_1 + \dots + X_n}{n}\right] = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \quad (\text{調査対象の分散/標本数})$$

2.3 標本平均の分布の意味

- ・ 標本平均で母平均を推定するのは平均的には正しい（バイアスがない，不偏である）
- ・ 標本数が増えるとそれに反比例して標本平均のばらつき = 分散が小さくなる

標本数が大きくなると標本平均は調査対象の平均にどんどん近づく(一致性がある)

2.4 標本数が十分大きいときの標本平均の分布

$$\text{中心極限定理から } \bar{X} = \frac{X_1 + \dots + X_n}{n} \longrightarrow \frac{\sigma}{\sqrt{n}}Z + \mu \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$(\sqrt{n}/\sigma)(\bar{X} - \mu)$ の分布が $Z \sim N(0,1)$, つまり標準正規分布に近づく .

3. 母平均の推定 (標本数が大きい場合)

95%信頼区間 = その区間の中に母平均が入る確率が95%になるような区間は,

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \text{ とすると, } \bar{X} - 1.96 \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{X} + 1.96 \frac{\hat{\sigma}}{\sqrt{n}}$$

一般に $(100 - \alpha)\%$ 信頼区間は, z_β を標準正規分布の 100 %点 (すなわち,

$$P(Z \leq z_\beta) = \beta \text{ を満たす } z_\beta) \text{ とすると } \bar{X} - z_{1-(\alpha/200)} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{1-(\alpha/200)} \frac{\hat{\sigma}}{\sqrt{n}}$$

比率の推定

n を標本数, m をあるカテゴリーに該当すると答えた標本の数とする. このとき, 調査対象中あるカテゴリーに属する対象の比率の95%信頼区間は,

$$\frac{m}{n} - 1.96 \sqrt{\frac{m}{n} \left(1 - \frac{m}{n}\right)} / \sqrt{n} \leq \mu \leq \frac{m}{n} + 1.96 \sqrt{\frac{m}{n} \left(1 - \frac{m}{n}\right)} / \sqrt{n}$$

(あるカテゴリーとは, 「大統領を支持する」, 「紅白を視聴した」とかで, その結果, 大統領支持率や紅白の視聴率がランダムサンプリング調査の結果から推定できる)

4. 母平均の推定 (正規分布の標本の場合)

$$95\% \text{ 信頼区間 } \bar{X} - t_{n-1, 0.975} S / \sqrt{n} \leq \mu \leq \bar{X} + t_{n-1, 0.975} S / \sqrt{n}$$

(ただし $S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$, $t_{n-1, \alpha}$ は自由度 $n - 1$ の t 分布の 100 %点)

5. 母分散の推定 (正規分布の標本の場合)

$$95\% \text{ 信頼区間 } \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\chi_{n-1, 0.975}^2} \leq \sigma^2 \leq \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\chi_{n-1, 0.025}^2}$$

$\chi_{n-1, \alpha}^2$ は自由度 $n - 1$ のカイ二乗分布の 100 % 点

スライドシリーズ7 仮説検定

2. 仮説検定の考え方

2.1 肯定したい仮説・否定したい仮説

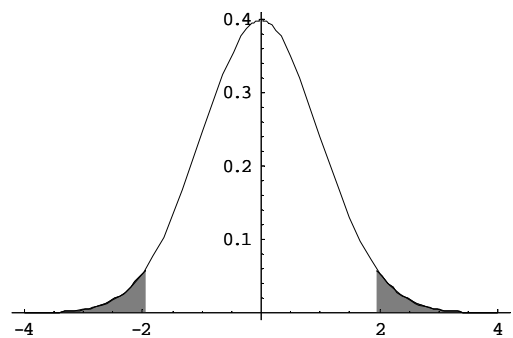
- ・ データを元に肯定的に実証したい仮説
対立仮説 (H_1) と呼ぶ
- ・ 「この仮説の否定」を「否定」することで、仮説を実証する
この否定されるべき仮説を帰無仮説 (H_0) と呼ぶ
- ・ 背理法に似た考え方

2.2 帰無仮説の棄却・受容 (1)

- ・ 検定の流れ
帰無仮説を「否定」する = 帰無仮説を棄却
棄却できれば、対立仮説を採択
= 対立仮説を「肯定」
棄却できなければ、帰無仮説を受容
= 帰無仮説を「肯定」
- ・ 棄却の考え方
推定値が < 帰無仮説が正しいとした場合に > 「確率論的に考えて起こりづらい値」になっていると棄却
右図の灰色の部分 = 棄却

2.3 棄却域と有意水準, 境界値

- ・ 灰色の範囲を棄却域と呼ぶ
- ・ 棄却域の決め方
灰色の部分の総面積を最初に決める
例えば 5% = 0.05 . 帰無仮説が正しい場合に帰無仮説を棄却する確率
この面積を有意水準と呼ぶ
5%, 1%, 10% などの順でよく使われる.
- ・ 棄却域との境目を境界値, 臨界値と呼ぶ



2.4 検定統計量

- ・ 影響されない分布を持つ形に変形し, 帰無仮説でのパラメータの値を代入したもの
この分布を元に判断

2.5 まとめ

- 検定統計量の式に、帰無仮説の値を全部代入する。

$$\text{帰無仮説： } \mu = 1 \text{ で } n=10 \text{ なら， } (\sqrt{10}/\hat{\sigma})(\bar{X}-1)$$

- 推定値を代入

$$\bar{X} = 1.31, \hat{\sigma} = 1.02 \text{ なら } (\sqrt{10}/1.02)(1.31-1) = 0.98$$

- この推定値が棄却域に入っているかを判断

有意水準 5% とすれば境界値は 1.96 だからその内側なので帰無仮説は受容される

$P(Z \leq 0.98) = 0.837$ なのでもし統計量値を境界値にしたら棄却域の確率は 0.326, 有意水準を上回るので受容

2.6 P 値

- P 値

統計量値を境界値にしたときの有意水準

これが実際の有意水準より大きければ、帰無仮説は受容

小さければ帰無仮説は棄却、対立仮説を採択

3. 母平均に関する検定

3.1 標本数が大きい場合

- 帰無仮説 $\mu = \mu_0$

- 検定統計量とその分布 $(\sqrt{n}/\hat{\sigma})(\bar{X} - \mu_0) \cong Z \sim N(0,1)$

3.1.1 両側検定 (対立仮説 型)

- 対立仮説が 型の時: $H_1: \mu \neq \mu_0$

推定値が仮説の値に対して大きい方に隔たっている場合、小さい方に隔たっている場合、帰無仮説を棄却する

棄却域は両側

- 境界値

有意水準を とすると, $\pm z_{1-\alpha/2}$

3.1.2 片側検定 (対立仮説 > 型)

- 対立仮説が > 型の時: $H_1: \mu > \mu_0$

推定値 \bar{X} が仮説の値に対して大きい方に隔たっていた場合だけ、帰無仮説を棄却する。

統計量が正の側が棄却域

棄却域は片側

- 境界値

有意水準を とすると $z_{1-\alpha}$

3.1.2 片側検定 (対立仮説 < 型)

・ 対立仮説が < 型 のとき: $H_1: \mu < \mu_0$

推定値 \bar{X} が仮説の値に対して小さい方に隔っていた場合だけ, 帰無仮説を棄却する

統計量が負の側が棄却域

棄却域は片側

・ 境界値

有意水準を α とすると, $-z_{1-\alpha}$