データ分析基礎 確率・統計の基礎

京都大学 国際高等教育院 附属データ科学イノベーション教育研究センター

世書といると関ア格人

sekido.hiroto.7a@kyoto-u.ac.jp

統計と確率の基礎知識

統計とは

- ★ 統計学とは、大雑把に言って、データの扱い方を考える学問
 - ★ データをどうやって得るか
 - ★ データをどうやって解析するか
 - ★ データを使って何か主張できるか
- ★ なんとなく、ではなく、できるだけ数学的(客観的)に議論する
- ★ 確率論を道具として用いることが多い
 - ★ データには誤差が付きもの(データを取る対象の選び方・測定誤差など)
 - ★ 正しいモデルは不明. 適当な近似を行いモデル化(モデル化誤差)
 - ★ 誤差の要因,振る舞いは複雑怪奇 → 確率的なものとして扱う
 - ★知りたいことにフォーカスを当て、わからないことは確率で暈す

確率とは

- ★○○を行ったとき、△△が起きる確率
 - **★**○○を**試行**,△△を**事象**と呼ぶ
- ★ 確率は、試行を行ったとき、その事象の「起こりやすさ」を実数で表したもの
 - $\star P($ 事象)やP(事象aが起こる)のように表す
 - ★確率は0以上1以下の実数で値が大きいほど起こりやすい
 - ★ 確率がpとは、十分大きなNに対して、試行を独立にN回行ったとき、だいたいNp回ぐらいその事象が起こるだろうと期待されるということ

確率とは

- ★確率として捉えるには
 - ★6月28日に東京で雨になる確率は53%である
 - ★ 気象データが残っている年を無作為に1つ選ぶとその年の6月28日の東京の天気が雨である確率が0.53である
 - ★日本人成人男性のうち身長が190cm以上の人の割合は0.06%である
 - ★日本人成人男性を無作為に1人選ぶと、その人の身長が190cm以上である確率が0.06%である
 - ★日本人男性が最近の環境で育つと、その人の身長が190cm以上となる確率が0.06%である
 - ★明日京都で雨が降る確率は10%である
 - ★ 観測しうる範囲で、今日と同じ気候条件であれば、その次の日京都で雨が降る確率が0.1である

条件付き確率

- ★ 試行を行い、事象 B が起こったとき、事象 A が起こる確率 P(A|B) または $P_B(A)$
 - $\star A$ も B も起こる確率 ÷ B が起こる確率, $P(A \cap B)/P(B)$ が定義
- ★ 試行にある種の制限をかけていると思うこともできる
 - ★今日と同じような気象条件になったとき、次の日が晴れる確率
 - ★ランダムに1日選び、選んだ日が今日と同じような気象条件であったときに、その次の日が 晴れである確率
- ★「確率」を考えるときは、(あるならば)得られている全ての情報を用い、条件付き確率を考える

余談1:モンティ・ホール問題

- ★ クイズ番組の優勝者が賞品を選ぶゲーム
- ★ドアが3つあり、開けるとそのうちの2つにはヤギ(外れ)、1つには車
- ★以下の手順を行う
 - ★優勝者は1つのドアを無作為に選ぶ
 - ★司会者はどれが外れか知っており、優勝者が選ばなかった外れのドアのうち1つを無作為に開く
 - ★優勝者は選ぶドアを変えることができる
 - ★ 最終的に優勝者が選んだドアに車があればもらえる
 - ★優勝者はどうすれば良いか? 車が貰える確率は?

余談1:モンティ・ホール問題

- ★ クイズ番組の優勝者が賞品を選ぶゲーム
- ★ ドアが3つあり、開けるとそのうちの2つにはヤギ(外れ)、1つには車
- ★以下の手順を行う
 - ★優勝者は1つのドアを無作為に選ぶ
 - ★司会者はどれが外れか知っており、優勝者が選ばなかった外れのドアのうち1つを無作為に開く
 - ★優勝者は選ぶドアを変えることができる
 - ★ 最終的に優勝者が選んだドアに車があればもらえる
 - ★優勝者はどうすれば良いか? 車が貰える確率は?
- ★ 適切な仮定のもとで…
 - ★選ぶドアを変えなければ車の確率1/3
 - ★選ぶドアを変えれば車の確率2/3

余談2:2つの封筒問題

- ★2つの封筒があり、中にはお金が入っている
- ★ 片方の封筒の中には片方の封筒の倍額入っているがどちらかわからない
- ★以下の手順を行う
 - ★1つの封筒を無作為に選び、中に何円入っているか確認する
 - ★選ぶ封筒を変更することができる
 - ★ 最終的に選んだ封筒の中身がもらえる
- ★無作為に選んだから最初選ばなかった方に倍額入っている確率は1/2
 - \star 選ぶ封筒を変えなければ確認した額x円もらえる
 - ★ 選ぶ封筒を変えれば貰える額の期待値は $\frac{1}{2}$ \frac
 - ★何がおかしいか?

余談3:為替?

- ★ 日本円をx円持っていて、1ドルx円でドルを買った
- \star 1ドルが2x円になるか、x/2円になったら1ドルを売る
- ★ 為替はランダムに動くとすると
 - ★ 円に対してドルの価値が2倍になることも、ドルに対して円の価値が2倍になることも、どちらが先に起こるかは同じ確率
 - \star 1ドルが2x円になるか、1ドルがx/2円になるか、どちらが先に起こるかは同じ確率
- ★ 最終的に手元に残る日本円の期待値は $\frac{1x}{22} + \frac{1}{2}(2x) = 1.25x > x$
 - ★何かおかしいのか?

母集団と標本

★ 母集団

- ★知りたい調査対象全体
- ★(確率)分布で表される

★ 標本

- ★母集団から抽出された集団
- $\star n$ 個の実数(の組)で表される
 - $\star n$ は標本サイズ

確率変数とは

- ★ 確率変数とは、試行を行うと値が定まるもの
 - ★サイコロを振ったときに出る目
 - ★全日本人の中から無作為に1人選んだとき、その人の身長
- ★ どんな値を取りうるか、またそれぞれの値はどれぐらい起こりやすいか、を両方表している
 - ★ 確率分布

確率変数とは

- ★ 取る値が連続的なものと離散的なものがある
 - ★ サイコロの出目は1,2,3,4,5,6のどれかなので離散的(飛び飛びの値)
 - ★身長は「この値またはこの値または…(可算無限)」と言えないので連続的と見なすことが 多い
 - ★ 少なくても連続的であると仮定して議論したほうが便利
 - \star ちょうど値がxになる確率,ということを考えないものは連続的
- ★ 複数の確率変数を考えても1つの試行
 - \star 全日本人の中から無作為に1人選んだとき,その人の身長Hと体重W
 - ★この場合、身長と体重は無関係ではない
 - ★ 身長が 170 cm のときに体重が 60 kg 以上である確率 $P(W \ge 60|H = 170)$, などと条件付き確率を考えることもできる

確率密度関数,確率質量関数

 \star 連続的な確率変数 X は確率密度関数 f(x) で記述

$$\star \int_{-\infty}^{\infty} f(x) \, \mathrm{d}x = 1, \quad f(x) \ge 0$$

 \star 離散的な確率変数 X は確率質量関数 f(x) で記述

$$\star \sum_{x} f(x) = 1, \quad f(x) \ge 0$$

$$\star P(X = a) = f(a) : f(a)$$
 はその確率変数が a となる確率

★2つ以上の確率変数を同時に考えれば、多変数関数になる

$$\star \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \, dx \, dy = 1$$

離散的な確率変数と確率質量関数の例

- ★(例) 理想的なサイコロの出目の場合
 - * 離散分布で、f(1) = f(2) = f(3) = f(4) = f(5) = f(6) = 1/6
 - ★ サイコロを振って出る目 X が 4 である確率は P(X=4)=f(4)=1/6. これは,何回もサイコロを振ると,その振った回数の 1/6 倍ぐらいの回数だけ 4 の目が出るという意味である.
- ★ 注意:通常サイコロは厳密に各目が出る確率が1/6となるわけではないと思われるが、1/6に非常に近いことが経験的に知られており、厳密にやるのは面倒、または、不可能なのでこういう風にしましょう、と考えている。

連続的な確率変数と確率分布関数の例

- ★(例)無作為に20歳の男性を1人選んだとき、その人の体重(kg)
 - ★ 連続分布で、f(60)とかf(70)ぐらいが大きくて、そこから外れると小さくなるであろう
- ★ 注意: 20歳の男性の人数は有限であるので、これは本当は離散分布であると思われる。しかし、厳密に体重を量れるわけでもなく、そもそも対象がかなり曖昧(今この瞬間の体重分布か? それとも人間が通常の生活を 20 年行った後の体重の分布か?)また、実際に必要となる値は、体重が 60kg 台の人口比率 $P(60 \le X < 70)$ などが多く、連続分布と考えたほうが便利な場合が多い。
 - ★実際の現象をできるだけうまく説明できる、数学的に扱うのが便利、などを考えて、うまく 設定する必要がある。

期待值、平均、分散、標準偏差

★ 確率 [密度 | 質量] 関数 f(x) に対応する確率変数 X の関数 s(X) の期待値

$$\star$$
 連続分布の場合: $\mathrm{E}[s(X)] = \int_{-\infty}^{\infty} s(x) f(x) \, \mathrm{d}x$

$$\star$$
 離散分布の場合: $\mathrm{E}[s(X)] = \sum_{x} s(x) f(x)$

★平均(mean, average),分散(variance),標準偏差(standard deviation)

- ★平均:E[X]
- ★ 分散: $V[X] = E[(X E[X])^2] = E[X^2] E[X]^2$

★標準偏差:
$$\sqrt{V[X]} = \sqrt{E[(X - E[X])^2]} = \sqrt{E[X^2] - E[X]^2}$$

★ 分散の項の等式は $E[\alpha s(X) + \beta t(X)] = \alpha E[s(X)] + \beta E[t(X)] \Rightarrow E[1] = 1$ などから導かれる

期待値,平均,分散,標準偏差

- ★ 平均,分散などは、その確率分布を知るための手がかりとなることがある
- ★ もしくは、確率分布関数を直接的に扱うのが難しいが、平均や分散は単なる実数なので扱いや すい
- **★ 平均**:確率変数がだいたいどのぐらいの値になるかという目安の1つ。 X_1, X_2, \ldots, X_N が独立で,それぞれその分布に従うなら, $(X_1 + X_2 + \ldots + X_N)/N$ は N が大きいとき,平均に近づく
- ★ 分散: 確率変数が平均からどれぐらい離れた値になるか、という目安の1つ。(平均からの距離の2乗)の平均、であるので、オーダーは「距離」の2乗
- ★ 標準偏差:分散の正の平方根. オーダーは「距離」

母集団 (population) と標本 (random sample)

- $\bigstar X_1, X_2, \ldots, X_N$ が独立で、すべて同じ確率分布に従うとき、これをサイズNの標本という
 - ★独立であるとは、荒っぽく言うと、 X_1 がとある値のとき、 X_2 はとある値になりやすい、などといった関係がないという事。標本が従う確率分布を母集団分布という
- \bigstar (例)サイコロをN回振るとき,i回目に出た目の数を X_i で表すと, X_1, X_2, \ldots, X_N はサイズNの標本となる
 - ★1回目のサイコロの出目は、2回目のサイコロの出目に影響しないであろうから、これらは独立である
- **★注意**:統計処理を行う際には、標本 X_i には、与えられたデータなどの実数が入っている。ただし、数学的に、標本や、その標本から得られる量はどのような性質があるかなどを議論したいときは、標本を確率変数と考えている

標本平均,標本分散,不偏分散

★ サイズ N の標本 X_1, X_2, \ldots, X_N に対して,以下が定義される

$$\bigstar$$
標本平均: $\overline{X} = \frac{X_1 + X_2 + \cdots + X_N}{N}$

★標本分散:
$$\frac{1}{N} \sum_{i=1}^{N} (\overline{X} - X_i)^2 = \overline{X^2} - (\overline{X})^2$$

$$\star$$
 不偏分散: $\frac{1}{N-1}\sum_{i=1}^{N}(\overline{X}-X_i)^2$

★ 注意:標本平均の期待値は母集団分布の平均に一致する。また、不偏分散の期待値は、母集団分布の分散に一致する。標本に対して分散を調べるときは標本分散を、標本を通じて、母集団分布の分散を調べるときは不偏分散を用いることが多い。

他の代表値(確率変数の大体の大きさを表す)

★ 中央値(median)

- ★標本(データ)に対して、小さい順に並び替えたときに真ん中の値(真ん中が2つあるなら、 足して2で割る)
- \star 確率変数 X に対して, $P(X \leq c) = 0.5$ なる c
 - ★連続分布に対して,

$$\int_{-\infty}^{c} f(x) \, \mathrm{d}x = 0.5$$

なるc

★ 確率変数 X に対して, $P(X \le c) = \alpha$ なる c を上側 α 点,もしくは, 100α % 点という

★ 最頻値(mode)

- \star 標本に対して、 X_1, X_2, \ldots, X_n のうち、最も多く登場する値
 - ★実際には、度数分布において、最も該当する数が多い階級とすることが多い
- \star 連続分布、離散分布に対して、f(x)の最大値を取るx