

数理統計 補助資料 ～ベイズ統計学とは？～

2023年度2学期： 月曜1限, 水曜3限
担当教員： 石垣 司

1

ベイズ統計学とは？

ベイズの定理を積極的に利用する統計学

- 人間の主観や信念などを情報創出のために利用できる

この授業で紹介するベイズ統計学の有用性

1. 統計的意思決定への応用
2. 情報の逐次更新
3. 柔軟な統計モデル

どれもが深奥な学術的意義と応用上の有用性をもつ。
本授業ではベイズ統計学の核となる部分を紹介



#1 非ベイズ統計学を頻度論や頻度主義的統計学とよぶ。信頼区間、統計的検定などは頻度論的な方法論
#2 頻度主義的統計学とベイズ統計学は哲学的なアプローチが異なる(端的には“確率”や“科学的”の考え方が違う)ため、長い間論争が絶えなかった。ただ、最近ではベイズ統計学の有用性が周知してきたため論争は落ち着いている。どちらが正しい・優れているという話ではない。考え方、有用な場面が異なる。

2

補足:ベイズ統計学に関する本授業の内容

ベイズ統計学の核となる部分と有用性の紹介

- ベイズ推定, 特に事後分布の推定に関する数理的な内容は学部2年生の学習内容を大きく超えてしまうため詳細の説明は行うことができない。興味のある学生はガイダンスで紹介した参考書を参考にしてほしい。
- とはいえ, How to 本ではないので, 本質的な部分の説明は数学も含めてしっかり紹介する。

3

(この授業で扱う)統計的な因果関係 #1

統計的因果関係(この授業で扱う因果関係)の表現

- ある原因に応じて結果が確率的に出現する因果関係

原因と結果の事象

結果となる事象の集合 $A = \{a_1, \dots, a_{|A|}\}$

結果の原因となる事象の集合 $B = \{b_1, \dots, b_{|B|}\}$

$a_1, \dots, a_{|A|}$ は互いに排反。 $b_1, \dots, b_{|B|}$ は互いに排反



【原因】

- インフルエンザワクチン接種(する, しない)
- 居住地300km以内に低気圧(あり, なし)
- 意中の人へ告白(する, しない)

【結果～不確実性をもつ】

- ⇒ インフルエンザ発症 or 非発症
- ⇒ 頭痛が起こる or 起こらない
- ⇒ 恋人になる or ならない(保留等含む)



#1 因果関係は必然関係(BならばAでなければならぬ)ではないので注意
#2 因果関係の定義、統計的因果推定、因果効果は本セメスターの終わりの回で取り扱う予定

4

原因に対する結果の出やすさを確率で表現する


- 例: 病気と検査結果

$A = \{a_1, a_2\}$ (a_1 :検査で陽性, a_2 :検査で陰性)

$B = \{b_1, b_2\}$ (b_1 :罹患している, b_2 :罹患していない)

$\Pr(B = b_1)$: ある人が罹患している確率 or 集団の中で罹患している人の割合など、分析の文脈で使い分ける

$\Pr(A = a_1|B = b_1)$: 罹患している人が検査で陽性になる確率



		Pr(A B)				
		B = b ₁	B = b ₂			
A = a ₁		0.9	0.2	Pr(B)	B = b ₁	0.4
		0.1	0.8		B = b ₂	0.6

重要な注意事項

- 数式や確率表自体は原因と結果を区別&規定できない
- 確率や条件付確率の意味は内容に合わせて分析者が解釈
「ワクチンを打ったけどインフル発症」or「ワクチンを打ったからインフル発症」



Thomas Bayes 1702-1761

ベイズの定理

- 事後確率(逆確率)を求める公式

$$\Pr(B|A) = \frac{\Pr(A|B)\Pr(B)}{\Pr(A)} = \frac{\Pr(A|B)\Pr(B)}{\sum_{i=1}^{|B|} \Pr(A|B = b_i)\Pr(B = b_i)} \quad \text{check!}$$

- $\Pr(B|A)$: 事後確率, $\Pr(B)$: 事前確率

事後確率 $\Pr(b_1|a_1)$ の解釈(今後, $\Pr(B = b_1|A = a_1)$ を $\Pr(b_1|a_1)$ と記述)

- ある結果 a_1 が生じる原因として互いに排反な b_1 or b_2 のみが考えられる場合、結果 a_1 が生じた原因が b_1 である確率 $\Pr(\text{原因})$, $\Pr(\text{結果}|\text{原因})$ が既知 $\Rightarrow \Pr(\text{原因}|\text{結果})$ を計算可能
- $\Pr(b_1|a_1)$: 検査で陽性だった人が罹患している確率

ベイズの定理を用いた例題 #1

例題: ある検査薬を用いた病気の検査を考える

- 罹患者が陽性反応を示す確率 99.9%
- 非罹患者が陽性反応を示す確率 0.2%
- その病気の日本人の罹患率は100,000人に1人
- 日本人全体から無作為抽出したある人にその検査薬を用いたところ陽性反応が出た。その人が罹患している確率は?

- 原因: b_1 :罹患, b_2 :非罹患
- 結果: a_1 :陽性反応, a_2 :陰性反応
- 事前確率: $\Pr(b_1) = 0.00001$
- その他の確率: $\Pr(b_2) = 0.99999$, $\Pr(a_1|b_1) = 0.999$, $\Pr(a_1|b_2) = 0.002$

ベイズの定理を用いた例題 #2

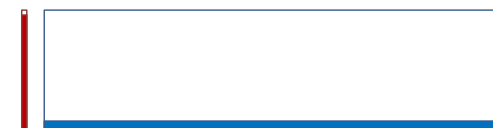
例題の解答(求めるべきは事後確率 $\Pr(b_1|a_1)$)

$$\Pr(b_1|a_1) = \frac{\Pr(a_1|b_1)\Pr(b_1)}{\Pr(a_1|b_1)\Pr(b_1) + \Pr(a_1|b_2)\Pr(b_2)} = \frac{0.999 \times 0.00001}{0.999 \times 0.00001 + 0.002 \times 0.99999} = \frac{999}{200997} \approx 0.005$$

- 陽性反応の出た人が罹患している確率は 0.5%

母集団: 日本人全体

罹患者: 0.001% 非罹患者: 99.999%
陽性反応率: 99.9% 陽性反応率: 0.2%



陽性反応が出た人の罹患率 = $\frac{\text{赤}}{\text{青}}$

#メモ 実用上は、「日本人全体から無作為抽出」と「症状が出たので病院に来た人」では統計的な問題も医療上の意義も異なるので注意

ベイズの定理による情報創出の構造

例題の解答から得られる示唆

- 事前確率の値が事後確率の値に影響を与える
- 事後確率をもつ情報は増えている
 検査を受けていない人の罹患率: 0.0001 %
 陽性反応が出た人の罹患率: 0.5 % \Rightarrow 危険度は検査前の5,000倍

ベイズの定理による情報創出の構造

事後の情報 \leftarrow データが持つ情報 + 事前の情報

- ベイズの定理は確率の公理を満たす単なる数式であるが、この構造を自然に表現している
 以降は、この情報創出の構造の統計モデルでの利用を考える

9

ベイズ統計学と統計モデル

統計モデル(復習)

- 関数 $f(y; \theta)$ と確率的構造を組み合わせる
- 関数の形はパラメータ θ の値で決まる
- 統計モデルの推定: データを用いて θ の値を推定

ベイズ統計学での統計モデル

- パラメータ θ を確率変数とみなして事前情報をモデル化する
- 事前情報+データを用いて θ の値を推定

例: 単回帰モデル $y = a + bx$

- もしも、パラメータ a と b について何らかの情報を持っているとき、 a と b を確率変数とみなして、確率密度関数でその情報を表現

10

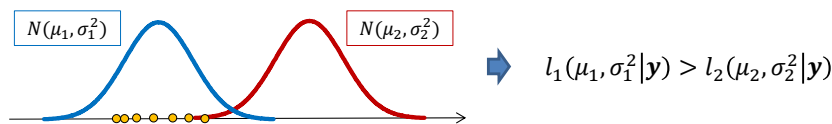
数学的準備 #1 (復習)

尤度: 与えられたデータがある確率分布から発生していると考えたときの尤もらしさの度合い

- データ y_i を発生させる確率分布: $p(y_i; \theta)$
 (パラメータ $\theta = [\theta_1, \dots, \theta_p]^T$)
- データ y_i の尤度: $l(\theta|y_i) = p(y_i; \theta)$
- 独立同一分布から発生したデータ $y = [y_1, \dots, y_N]^T$ の尤度:

$$l(\theta|y) = \prod_{i=1}^N p(y_i; \theta)$$

例: 次のデータは青と赤のどちらの正規分布から発生していると考えるのが妥当か?



11

数学的準備 #2

データとパラメータの確率分布の表記方法

- $p(y)$: データベクトル $x = [y_1, \dots, y_N]^T$ が生じる同時確率密度関数。 $p(y) = p(y_1, \dots, y_N)$
- $p(\theta)$: パラメータベクトル $\theta = [\theta_1, \dots, \theta_p]^T$ の各要素を確率変数とみなしたときに θ が生じる同時確率密度関数。 $p(\theta) = p(\theta_1, \dots, \theta_p)$
 例: 正規分布では $p(\theta) = p(\mu, \sigma^2)$
- 確率の乗法定理: $p(y, \theta) = p(y|\theta)p(\theta)$
- $p(y) = \int p(y, \theta) d\theta = \int p(y|\theta)p(\theta) d\theta$
- $\int p(y|\theta)p(\theta) d\theta$: パラメータの集合で積分をする記号
 例: $\theta_i \in \mathbb{R} (i = 1, \dots, P)$ なら、
 $\int p(y|\theta)p(\theta) d\theta = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} p(y|\theta)p(\theta) d\theta_1 \dots d\theta_p$

12

データ&パラメータ&ベイズの定理 #1

データとパラメータを結びつけるベイズの定理

- 連続型確率変数のベイズの定理

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)} = \frac{p(y|\theta)p(\theta)}{\int p(y|\theta)p(\theta)d\theta}$$

- $p(\theta|y)$: 事後確率分布(事後分布)
- $p(\theta)$: 事前確率分布(事前分布)
- $p(y|\theta)$: 尤度関数
- $p(y)$: 周辺尤度

13

データ&パラメータ&ベイズの定理 #2

尤度の観点からみたベイズの定理の意味

- パラメータ θ を確率変数と見ている
- 独立同一分布から発生したデータ $y = [y_1, \dots, y_N]^T$ を考えると, $p(y|\theta)$ は尤度関数 $\prod_{i=1}^N p(y_i; \theta)$
- 事後分布 $p(\theta|y)$ はデータ観測後のパラメータの分布

$$p(\theta|y) = \frac{\prod_{i=1}^N p(y_i; \theta) p(\theta)}{\int \prod_{i=1}^N p(y_i; \theta) p(\theta) d\theta} \propto \prod_{i=1}^N p(y_i; \theta) p(\theta)$$

\propto : 比例するという記号

- ベイズの定理の構造

事後分布 \leftarrow データ(尤度関数) \times 事前分布

14

演習問題

問題

- ある企業では同じ型の製品を工場A,B,Cでそれぞれ生産し販売している。また、不良品出荷率が少ないほど生産シェアを高く設定している。
- 生産シェア 工場A:70%, 工場B:25%, 工場C:5%
- 不良品出荷率 工場A:0.1%, 工場B:0.2%, 工場C:0.4%

- ある1つの販売済みの製品が不良品と判明した。その不良品が工場Aでつくられた製品である確率を求めなさい。

15