

# ベイズ統計

古谷知之

## 講義概要

- 自然共役事前分布
- 一変量正規分布のベイズ推定
  - 正規分布の平均のベイズ推定
  - 正規分布の分散のベイズ推定
- 経験ベイズ推定と階層ベイズ推定
- 尤度原理

## 自然共役事前分布

- 事前分布と事後分布が、互いに似たような確率分布を持つような場合に、その事前分布を自然共役事前分布という
- 主な自然共役事前分布は以下のとおり

事前分布	尤度関数	事後分布
二項分布	ベータ分布	ベータ分布
ポアソン分布	ガンマ分布	ガンマ分布
正規分布(平均)	正規分布	正規分布
正規分布(分散)	正規分布	逆ガンマ分布

## 試合の得失点差のベイズ推定

- あなたなあるプロスポーツチームのアナリスト兼ストラテジストだとします。
- チームの攻撃力と守備力を示す統計量(stats)に試合の(ホームチームとアウェイチームの)得失点差(H-A)を用いています。
- 今シーズンのあなたのチームの攻撃力と守備力を試すために、プレシーズンマッチ5試合を戦いました。

NFL日本公式サイト  
NFL JAPAN.COM

SHOP RECOMMEND  
ドラフト時に選手が被るあのキャップが遂に解禁!

Sportsnavi

Metaphors  
Selected Writings of Michael Ray Fitzgerald  
A way of life – by Michael Fitzgerald

Hear Porter Robinson, Post-EDM's Greatest Hope, Duet with a Robot on

Arresting Children Is Now Commonplace in America United

スポーツナビ TOP ニュース コラム フォト 動画 テレビ放送 スケジュール **スコア** 順位表 成績 チーム 特集

About NFL ルール入門 日本人選手 チアリーダー イベント情報 ブログ ファンクラブ 国内アメフト More NFL ショップ

TOP > スコア

スコア > シーズン日程表

2014年

プレシーズン					レギュラーシーズン																
HOF	1週	2週	3週	4週	1週	2週	3週	4週	5週	6週	7週	8週	9週	10週	11週	12週	13週	14週	15週	16週	17週

ポストシーズン

ワイルドカード ディビジョナル AFC/NFCチャンピオンシップ プロボウル スーパーボウル

● レギュラーシーズン 第1週

\*試合開始時間は日本時間表示  
\*下がホームチーム

PR  
Saddle Up for Summer 2014

スコアボード テレビ放送

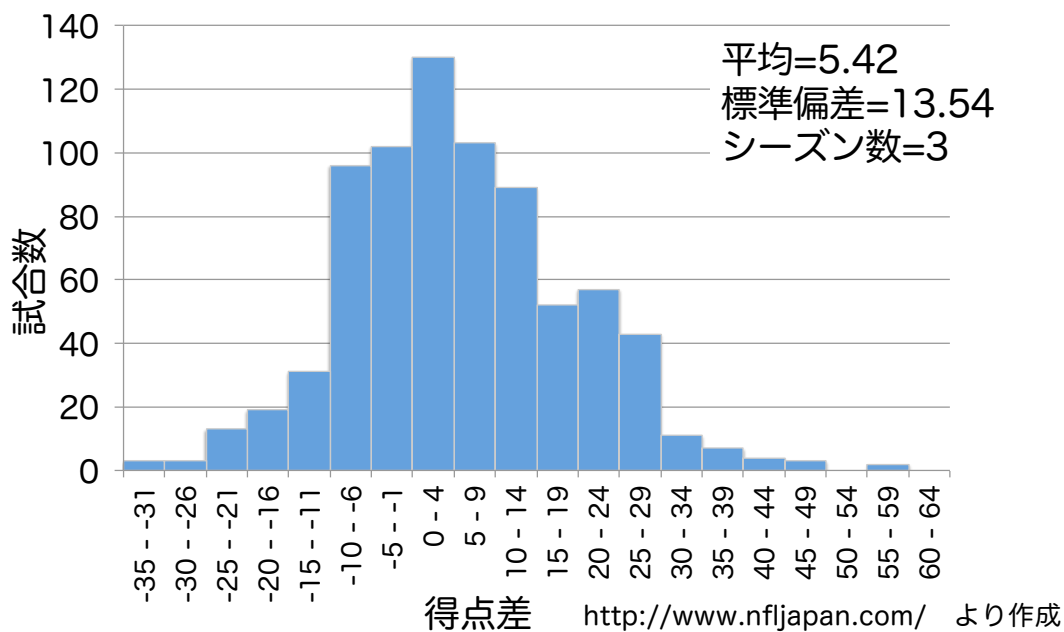
SHOP MOBILE

<http://www.nfljapan.com/>

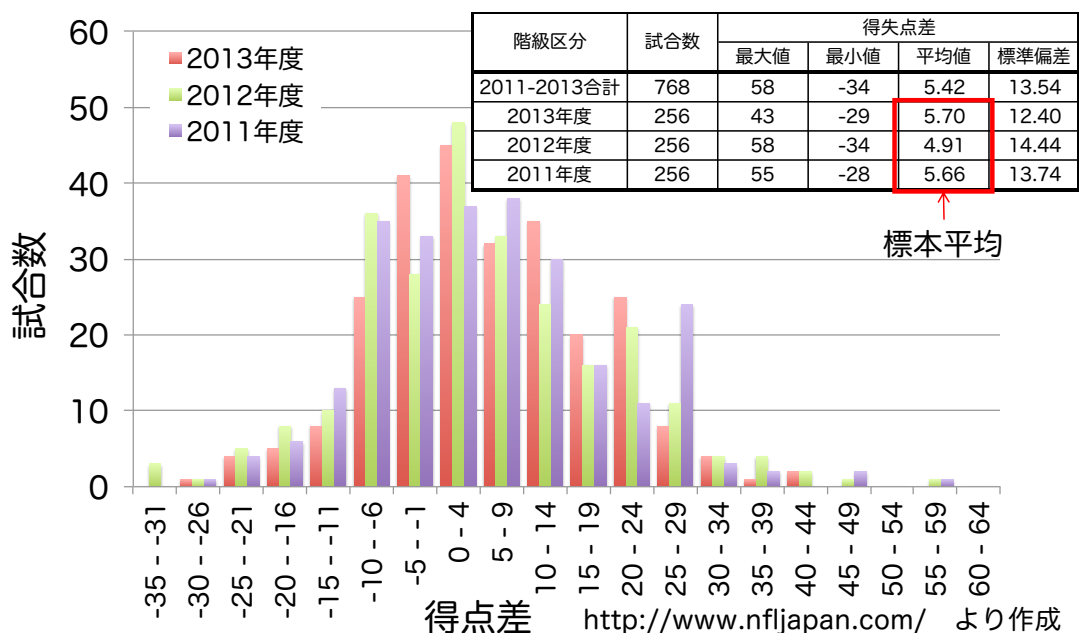
## 試合の得失点差のベイズ推定

- プレシーズンマッチでは、得失点差の平均と標準偏差は以下のとおりでした。
  - 平均 = 6.00、標準偏差 = 13.62、5試合分
- 他方、過去3年間のリーグ全体の得失点差の平均と標準偏差は以下のとおりでした。
  - 平均 = 5.42、標準偏差 = 13.54、768試合分
- あなたのチームの得点差の平均について、少ないプレシーズンマッチの結果から何か言えるのだろうか？

# NFL2011-2013シーズンの 試合別得点差分布(H-A)



# NFL2011-2013シーズンの 試合別得点差分布(H-A)



## 試合の得失点差のベイズ推定

- 得失点差の分布はある値（平均値）を中心に左右対称の釣鐘型の分布をしているように見えます。
- このような形をする確率分布に「正規分布」があります。

## 正規分布（ガウス分布）

- 平均 $\mu$ 、分散 $\sigma^2$ となる確率変数 $x$ について、以下の確率密度関数に従う分布を正規分布 $N(\mu, \sigma^2)$ という

$$N(\mu, \sigma^2) = f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- 正規分布の平均と分散は以下のとおり

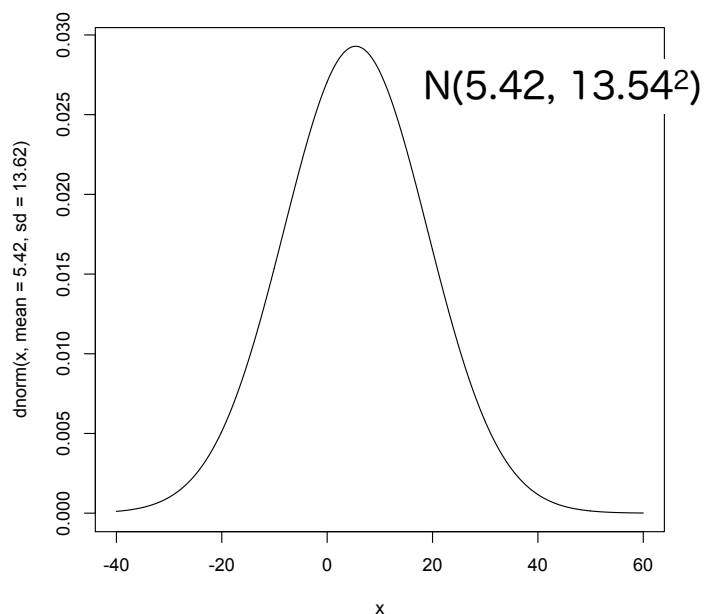
$$\begin{aligned} E(X) &= \mu \\ V(X) &= \sigma^2 \end{aligned}$$

## 得失点差の確率分布

- 2011-2013年の3シーズン(全768試合)のデータから、平均5.42、標準偏差13.54の正規分布 $N(5.42, 13.54^2)$ は

$$N(5.42, 13.54^2) = \frac{1}{\sqrt{2\pi \times 13.54^2}} \exp\left(-\frac{(x - 5.42)^2}{2 \times 13.54^2}\right)$$

## 得失点差の確率分布



## 試合の得失点差のベイズ推定

- 得失点差に関する事前分布と尤度関数を正規分布で与える
  - 事前分布：プレシーズンマッチ5試合分
  - 尤度関数：2011-2013年の3シーズン分
- 事後分布として得失点差の平均をベイズ推定したい
  - プレシーズンマッチと2011-2013シーズンのデータから、今シーズンの攻撃力・守備力を占いたい。
  - プレシーズンマッチの試合結果から判断するのは拙速ではないか？

## 試合の得失点差のベイズ推定

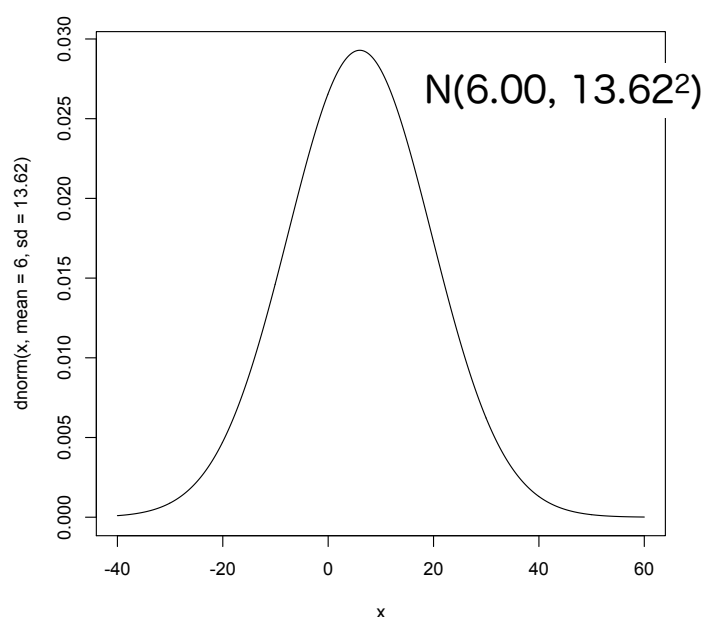
- 事前分布

$$\mu \sim N(6.00, 13.62^2)$$

$$\pi(\mu) = \frac{1}{\sqrt{2\pi \times 13.62^2}} \exp\left(-\frac{(\mu - 6.00)^2}{2 \times 13.62^2}\right)$$

$$\propto \exp\left(-\frac{(\mu - 6.00)^2}{2 \times 13.62^2}\right)$$

## 得失点差の確率分布



## 試合の得失点差のベイズ推定

- 尤度関数
- 過去の全試合(母集団)から3シーズン分を標本分布として切り出したデータを使う
- 各シーズンのデータを $X_i$  ( $i = 1, 2, 3$ )とする

$$X_1, X_2, X_3 \sim N(5.42, 13.54^2)$$

$$f(x_i|\mu) = \frac{1}{\sqrt{2\pi \times 13.54^2}} \exp\left(-\frac{(x_i - 5.42)^2}{2 \times 13.54^2}\right)$$



## 試合の得失点差のベイズ推定

- 尤度関数
- 3シーズン分の合計⇒正規分布の積

$$\begin{aligned}\prod_{i=1}^3 f(x_i|\mu, 13.54) &\propto \prod_{i=1}^3 \exp\left(-\frac{(x_i - \mu)^2}{2 \times 13.54^2}\right) \\ &= \exp\left(-\frac{\sum_{i=1}^3 (x_i - \mu)^2}{2 \times 13.54^2}\right)\end{aligned}$$

## 試合の得失点差のベイズ推定

- 事後分布

$$\begin{aligned}f(\mu|X) &\propto \exp\left(-\frac{\sum_{i=1}^3 (x_i - \mu)^2}{2 \times 13.54^2}\right) \cdot \exp\left(-\frac{(\mu - 6.00)^2}{2 \times 13.62^2}\right) \\ &\propto \exp\left(-\frac{\mu'}{2 \times \sigma'^2}\right)\end{aligned}$$

# 試合の得失点差のベイズ推定

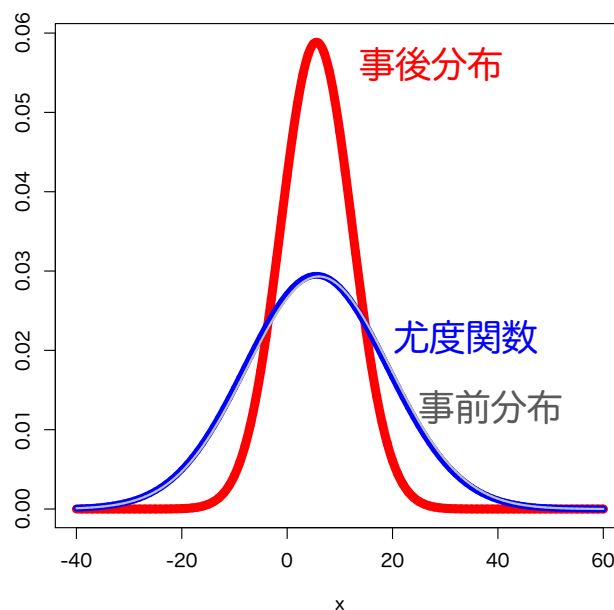
- 事後分布

$$\pi(\mu|X) \propto \exp\left(-\frac{\mu'}{2 \times \sigma'^2}\right)$$

$$\mu' = \frac{6.00/13.62^2 + 3 \cdot 5.42/13.54^2}{1/13.62^2 + 3/13.54^2} \approx 5.56$$

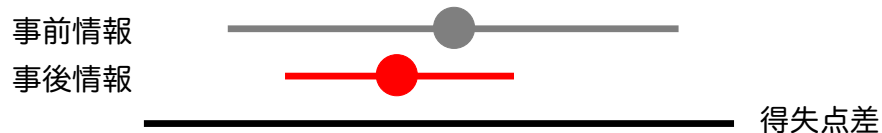
$$\frac{1}{\sigma'^2} = \frac{1}{13.62^2} + \frac{3}{13.54^2} = 6.77^2$$

# 試合の得失点差のベイズ推定



## 事前情報だけで判断しない理由

- 平均と分散(標準偏差)だけでなく、95%信頼区間も異なる
- 事後情報を採用するのがベイズアン



	平均	標準偏差	95%信頼区間
事前情報	6.00	13.62	[-20.7, 32.7]
事後情報	5.56	6.78	[-7.7, 18.9]

## 正規分布の平均のベイズ推定

- 事前分布(分散既知の正規分布)  $\mu \sim N(\mu_0, \sigma_0^2)$

$$\pi(\mu) \propto \exp\left(-\frac{(\mu - \mu_0)^2}{2 \times \sigma_0^2}\right)$$

- 尤度関数(標本分布が正規分布)  $X_i \sim N(\mu, \sigma^2)$

$$\prod_{i=1}^n f(x_i | \mu, \sigma^2) \propto \prod_{i=1}^n \exp\left(-\frac{(x_i - \mu)^2}{2 \times \sigma^2}\right) \propto \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2 \times \sigma^2}\right)$$

## 正規分布の平均のベイズ推定

- 事後分布

$$f(\mu|X) \propto \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2 \times \sigma^2}\right) \cdot \exp\left(-\frac{(\mu - \mu_0)^2}{2 \times \sigma_0^2}\right) \\ \propto \exp\left(-\frac{\mu_1}{2 \times \sigma_1^2}\right)$$

- ここで、

$$\mu_1 = \frac{\mu_0/\sigma_0^2 + n \cdot \bar{x}/\sigma^2}{1/\sigma_0^2 + n/\sigma^2}$$

$$\frac{1}{\sigma_1^2} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}$$

## 正規分布の平均のベイズ推定

- 事後分布導出の過程で

$$\begin{aligned} \sum_{i=1}^n (x_i - \mu)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \\ &= (n-1)s^2 + n(\bar{x} - \mu)^2 \\ &\quad \uparrow \\ &\quad \text{標本分散} \\ &\quad = \text{定数} \end{aligned}$$

# 正規分布の平均のベイズ推定

## クイズ

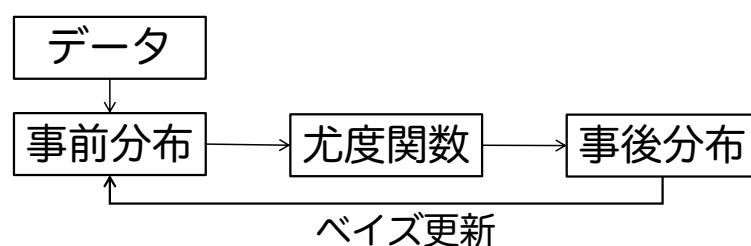
環境汚染が問題となっているある地点で、環境汚染物質の測定を5回の行ったところ、平均=3.2、分散=4.8<sup>2</sup>であった。この環境汚染物質は、同じ時点・場所・測定器で測定しても値がばらつくことが知られている。

これまで200箇所での測定を通じて、その分布が正規分布  $N(2.5, 4.5^2)$  に従うことが知られている。

上述の測定地点における環境汚染物質の事後分布平均値(事後平均)を求めなさい。

## 経験ベイズ推定

- 事前分布をデータから与える方法を経験ベイズ (empirical Bayes) 推定という
- 事前分布のパラメータを超パラメータ (hyper parameter) という
  - 事前分布の分散(精度)など



## 精度(precision)

- 正規分布をもちいてベイズ推定する際に、分散の逆数を精度として表すことがある

$$\tau = \frac{1}{\sigma^2}$$

- ベイズ統計では、精度がより重要である
- 分散が小さい $\Leftrightarrow$ 精度が高い
- 事後正規分布の精度は、事前正規分布の精度とデータの精度の和
- 事後正規分布の平均は、事前正規分布の精度とデータの精度の平均

## 正規分布の分散のベイズ推定

- 正規分布のベイズ推定における尤度関数

$$\begin{aligned} \prod_{i=1}^n f(x_i|\mu, \sigma^2) &\propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &\propto \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right) \end{aligned}$$

- ここで $\alpha - 1 = n/2$ ,  $\lambda = \sum_{i=1}^n (x_i - \mu)^2 / 2$ と置き換えると、

$$\prod_{i=1}^n f(x_i|\mu, \sigma^2) \propto \left(\frac{1}{\sigma^2}\right)^{\alpha-1} \exp\left(-\frac{\lambda}{2\sigma^2}\right)$$

## 正規分布の分散のベイズ推定

- さらに、 $\tau = \frac{1}{\sigma^2}$  と置き換えると

$$\prod_{i=1}^n f(x_i | \mu, \sigma^2) \propto (\tau)^{\alpha-1} \exp(-\lambda \cdot \tau)$$

- これはガンマ関数  $Ga(\alpha, \lambda)$  に比例する

$$Ga(\alpha, \lambda) \propto (\tau)^{\alpha-1} \exp(-\lambda \cdot \tau)$$

- すなわち事前分布の精度  $\tau$  に対して、自然共役事前分布としてガンマ分布を用いることと同じ意味を持つ

## ガンマ関数

- ガンマ関数  $\Gamma(\alpha)$  は次式で表される

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$$

$$\alpha > 0$$

- ガンマ関数  $\Gamma(\alpha)$  は以下のような性質を持つ

$$\begin{aligned}\Gamma(\alpha) &= (\alpha - 1)! \\ \Gamma(\alpha + 1) &= \alpha \Gamma(\alpha) \\ \Gamma(1/2) &= \sqrt{\pi}\end{aligned}$$

# ガンマ分布

- ガンマ分布の確率密度関数 $f(x)$ は次式のようになる

$$f(x) = Ga(\alpha, \lambda) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, x \geq 0$$

- ガンマ分布の確率変数を $X$ とすると、平均 $E(X)$ と分散 $V(X)$ は以下の通り

$$E(X) = \frac{\alpha}{\lambda}$$

$$V(X) = \frac{\alpha}{\lambda^2}$$

## 正規分布の分散のベイズ推定

- さらにいえば、事前分布の分散 $\sigma^2$ に対する自然共役事前分布として、逆ガンマ分布を用いることと同じ意味である。
- 逆ガンマ分布は、正規分布のベイズ推定にはよく用いられる
- もっとも、事前分布の分散の事前分布として逆ガンマ分布を用いず、単に分散の逆数などとして与える場合もある=無情報事前分布



# 逆ガンマ分布

- 逆ガンマ分布

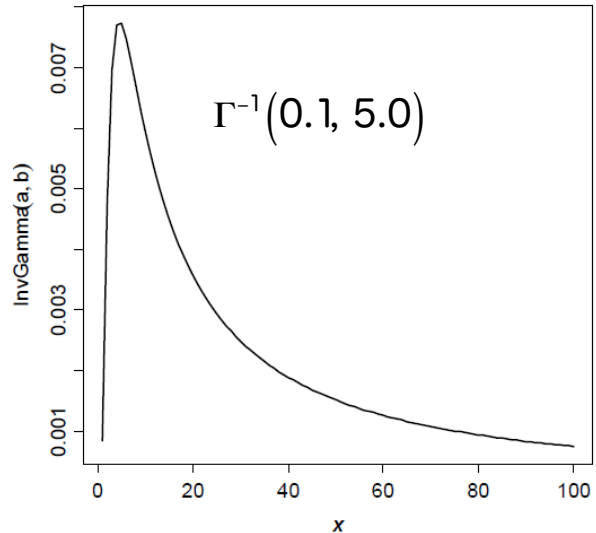
$$Ga^{-1}(\alpha, \lambda)$$

- 性質

- 平均： $\lambda/(\alpha - 1), \alpha > 1$
- 分散：

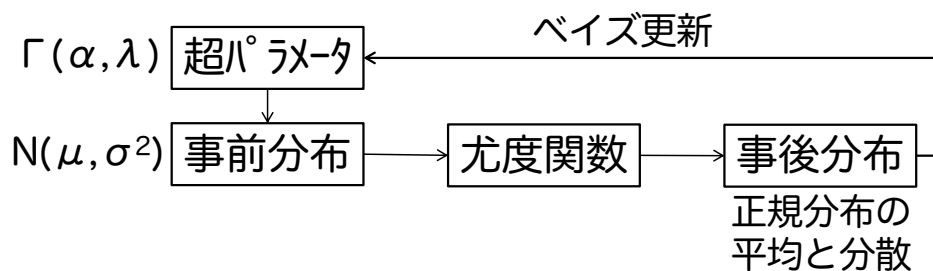
$$\frac{\lambda^2}{(\alpha - 1)^2(\lambda - 1)}$$

$\lambda > 1, \alpha > 2$



# 階層ベイズ推定

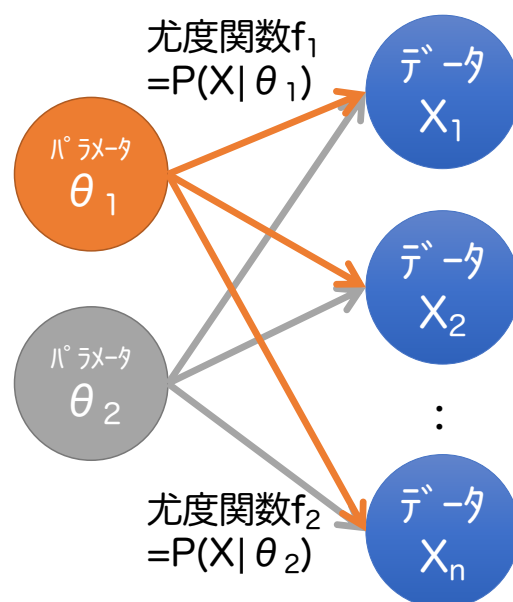
- 事前分布をベイズ推定する方法を階層ベイズ (hierarchical Bayes) 推定という
- 事前分布の超パラメータをベイズ推定する
  - 事前正規分布の分散(逆ガンマ分布に従う)の超パラメータのベイズ推定



## ベイズ推定を行う上での疑問

- 事前分布はどのような確率密度分布でもよいのか？
- ベイズ更新の際、事前分布はどの順番で与えても同じ事後分布が得られるのか？
- 同じ事前分布についてどのような尤度関数でも同じ事後分布が得られるのか？

尤度：パラメータからデータを再現する



## 尤度原理

- 未知パラメータ $\theta$ を推定するとき、尤度に観測されたデータの全ての情報を与えている場合には、比例関係にある尤度関数からは、同じ事前分布に対して同じ推論結果が導かれる
- ベイズ推定法は尤度原理を満たす
- 最尤推定法は尤度原理を満たさない

## 尤度原理

- コイントスをして10回中6回が表だった
- このとき尤度関数 $f_1$ は
$$f_1(x|\theta) = {}_{10}C_6 \theta^6 (1-\theta)^4 \propto \theta^6 (1-\theta)^4$$
- 9回コインを投げて5回表が出た後に10回目に表が出る確率密度関数 $f_2$ は
$$f_2(x|\theta) = \binom{10-1}{6-1} \theta^{6-1} (1-\theta)^4 \propto \theta^6 (1-\theta)^4$$
- この2つの関数は比例関係にある
$$f_1(x|\theta) \propto f_2(x|\theta)$$

## 尤度原理

- ちなみに頻度主義統計学においては、帰無仮説 $H_0$ と対立仮説 $H_1$ を次のようにおく

$$H_0: \theta = 0.5, H_1: \theta > 0.5$$

- 帰無仮説 $H_0$ は有意水準5%で棄却される

$$\begin{aligned} P(X \geq 6 | H_0) &= \sum_{x=6}^{10} {}_{10}C_x \cdot 0.5^x \cdot (1 - 0.5)^{10-x} \\ &= \frac{210 + 120 + 45 + 10 + 1}{2^{10}} \approx 0.377 > 0.5 \end{aligned}$$

## 同一性と可換性

- 同一性(identification)

- すべてのデータ $x$  ( $\forall x \in \Omega$ )に対して、パラメータ $\theta_a, \theta_b$ が与えられた条件の下で $P(x|\theta_a) = P(x|\theta_b)$ が成立するとき、パラメータは同一 $\theta_a \equiv \theta_b$ でなくてはならない。

- 可換性(exchangeability)

- $n$ 回のベルヌーイ試行において、成功する回数 $s$ が同じであれば、成功・失敗が生じる順序がどうであれ確率は同じである

## デ・フィネッティの定理

- 事前分布がどのような確率密度分布であっても、どの順番で事前分布が出現しても、十分長い期間ベイズの定理を適用すると、同じ事後分布に収束する
- 可換性の性質が示されたことにより尤度と事前情報の存在が示された