

# ベイズ推定

植野真臣

電気通信大学 大学院

情報理工学研究科

# 今後のスケジュール(予定)

- 4月7日 授業の概要とガイダンス
- 4月14日 ベイズの定理
- 4月21日 ベイズはどのように誕生したか？
- 4月28日 ベイズはコンピュータ、人工知能の父である！！
- 5月12日 アランチューリングとベイズ
- 5月19日 ビリーフとベイズ
- 5月26日 尤度と最尤推定(1)
- 6月2日 尤度と最尤推定(2)
- 6月9日 ベイズ推定と事前分布(1)
- 6月16日 ベイズ推定と事前分布 (2)
- 6月 23日 データサイエンス：ルービン因果推論
- 6月30日 テストのデータサイエンス
- 7月7日 階層ベイズとデータサイエンス
- 7月14日 ベイジアンネットワークと因果推論
- 7月28日 国際会議で休講
- 8月 4日 テストと総括

# 1. ベイズ原理

## 定義1 (事後分布)

$X = (X_1, \dots, X_m)$  が独立同一分布  $f(x|\theta)$  に従う  $m$  個の確率変数とする. 確率変数に対応した  $n$  個のデータ  $x = (x_1, \dots, x_n)$  が得られたとき,

$$p(\theta|x) = \frac{p(\theta) \prod_{i=1}^n f(x_i|\theta)}{\int_{\Theta} p(\theta) \prod_{i=1}^n f(x_i|\theta) d\theta}$$

を事後分布 (posterior distribution) と呼び、  
 $p(\theta)$  を事前分布 (prior distribution) と呼ぶ.

# 注意: ベイズでの $\Theta$ の扱い

尤度では、 $\Theta$ は確率変数ではない  
ベイズでは事前・事後分布が確率  
法則に従うのであれば、 $\Theta$ は確率  
変数となる

$$p(\theta|x) = \frac{p(\theta) \prod_{i=1}^n f(x_i|\theta)}{\int_{\Theta} p(\theta) \prod_{i=1}^n f(x_i|\theta) d\theta}$$

# ベイズの推定での利点

ベイズでは、厳密な確率推論がパラメータ推定にも適用できる。

# 事後分布最大化推定量

## 定義2 (MAP推定値)

データ $x$ を所与として, 以下の事後分布最大となるパラメータを求めるとき,

$$\hat{\theta} = \arg \max \{p(\theta|x) : \theta \in C\}$$

$\hat{\theta}$ をベイズ推定値 (Bayesian estimator) または, **事後分布最大化推定値** (maximum a posterior estimator, **MAP 推定値**) と呼ぶ.

# EAP 推定値

## 定義3 (EAP 推定値)

データ $x$ を所与として, 以下の事後分布によるパラメータの期待値を求めるとき,  $\hat{\theta} = E(\theta|x)$ を期待事後推定値 (expected a posterior estimator, EAP 推定値) と呼ぶ.

ベイズ推定値も強一致性をもつ.

# ベイズ推定の一貫性

定理1 (ベイズ推定の一貫性)

ベイズ推定において推定値 $\hat{\theta}$ が最尤推定値に漸近的に一貫する場合、真のパラメータ $\theta^*$ の強一貫推定値となる。

定理12 (ベイズ推定の推定値の分散)

事後確率密度関数 $p(\theta|x)$ が以下で直接求められる。

$$\text{Var}(\theta|x)$$

# 重要な性質

$$E(\theta) = E_x(E_\theta(\theta|x))$$

$$\text{Var}(\theta) = E(\text{Var}(\theta|x)) + \text{Var}(E(\theta|x))$$

= モデル分散の期待値 + 期待値の分散

## 2. 無情報事前分布

$$p(\theta|x) = \frac{p(\theta) \prod_{i=1}^n f(x_i|\theta)}{\int_{\Theta} p(\theta) \prod_{i=1}^n f(x_i|\theta) d\theta}$$

を求めるための事前分布 $p(\theta)$ の設定について、どのように設定するかが問題となる。通常、データを採取するまで、われわれはデータについての情報をもたない。そのために、 $p(\theta)$ は無知を表す分布でなくてはならない。このような無知を示す事前分布を**無情報事前分布** (non-informative prior distribution) と呼ぶ。

# 無情報事前分布(Jeffreys1961)

母数 $\theta$ について,  $\theta \in (-\infty, \infty)$  のみの情報があるとき, 事前分布は一様分布となる.

$$p(\theta) \propto \text{const}$$

$\int_{-\infty}^{\infty} p(\theta) \neq 1$ となり, 事前分布 $p(\theta)$ は確率の公理を満たさない. このような事前分布を**improper prior distribution**と呼ぶ.

# 無情報事前分布(Jeffreys1961)

母数 $\theta$ について,  $\theta \in (0, \infty)$ のみの情報があるとき,  $\theta$ の対数が一様であるような事前分布を考える.  
すなわち,  $p(\log \theta) \propto \text{const}$ であるから, 変数変換すれば,

$$p(\theta) \propto \frac{1}{\theta}$$

$\int_{-\infty}^{\infty} p(\theta) \neq 1$ となり, **improper prior distribution**.

# 注) 変数変換 $\phi \rightarrow \theta$

$$p(\phi) \rightarrow p(f(\phi))$$

$\theta = f(\phi)$  とすると

$$p(\theta) = p(\phi) \frac{\partial \phi}{\partial \theta} = p(f^{-1}(\theta)) \frac{\partial \phi}{\partial \theta}$$

今、 $p(\phi) = p(\log \theta) \propto \text{const}$  より

$$p(\theta) = p(\phi) \frac{\partial \phi}{\partial \theta} = \text{const} \times \frac{1}{\theta} \propto \frac{1}{\theta}$$

Proper prior: principle of stable estimation (Edwards et al. 1963)

例えば,  $\theta \in [a, b]$  であれば,  $p(\theta) = \frac{1}{b-a}$  となり,  $\int_{-\infty}^{\infty} p(\theta) = 1$  と確率の公理を満たす.

$\theta \in [a, b]$  では,  $p(\theta) = \text{const}$  であるが,  $\kappa = \theta^{10}$  としても, ジェフリーズのルールに従えば,  $p(\kappa) = \text{const}$  となつてほしい. しかし, 変数変換すれば, そのようにならないことがわかる.

# Jefferys prior (Box and Tiao 1973)

パラメータ変換を許容するパラメータ空間でエントロピーを最大にする事前分布は

$$p(\theta) \propto \sqrt{I(\theta)}$$

$I(\theta)$  はフィッシャー情報量を示す.

これが、ジェフリーズが提唱した母数の変換の不変性から導いた分布に一致するので、**ジェフリーズの前分布**と呼ばれる。

# 自然共役事前分布(最も一般的！！)

これまでの事前分布では、データを得る前の事前分布と事後分布は、分布の形状が変化する。しかし、データの有無にかかわらず、分布の形状は同一のほうが自然。そこで、事前分布と事後分布が同一の分布族に属するとき、その事前分布を自然共役事前分布(natural conjugate prior distribution)と呼ぶ。

# 自然共役事前分布によるベイズ推定例

## 例1 (二項分布)

$$f(x|\theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

コインを投げて  $n$  回中  $x$  回表が出たときの

確率  $\theta$  をベイズ推定しよう.

尤度関数は、 $\binom{n}{x} \theta^x (1 - \theta)^{n-x}$ であり、

二項分布の自然共役事前分布は、以下のベータ分布 ( $Beta(\alpha, \beta)$ ) である。

$$p(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

事後分布は、 $p(\theta | n, x, \alpha, \beta) =$   
 $\frac{\Gamma(n + \alpha + \beta)}{\Gamma(x + \alpha)\Gamma(n - x + \beta)} \theta^{x + \alpha - 1} (1 - \theta)^{n - x + \beta - 1}$

とやはりベータ分布となる。

対数をとって、以下の対数事後分布を最大化すればよい。

$$\begin{aligned} & \log p(\theta | n, x, \alpha, \beta) \\ &= \log \frac{\Gamma(n + \alpha + \beta)}{\Gamma(x + \alpha)\Gamma(n - x + \beta)} \\ & \quad + (x + \alpha - 1)\log \theta \\ & \quad + (n - x + \beta - 1)\log(1 - \theta) \end{aligned}$$

$$\frac{\partial \log p(\theta|n, x, \alpha, \beta)}{\partial \theta} = 0 \text{ のとき,}$$

対数事後分布は最大となるので,

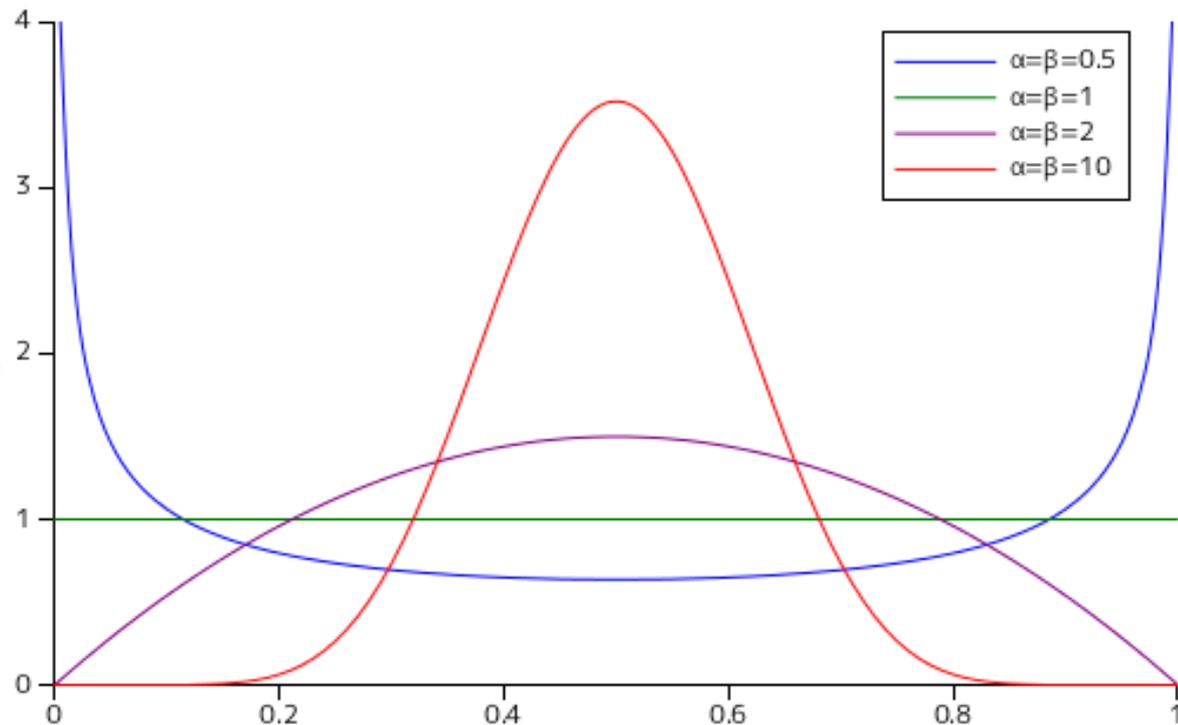
$$\begin{aligned} \frac{\partial \log p(\theta|n, x, \alpha, \beta)}{\partial \theta} &= \frac{(x + \alpha - 1)}{\theta} - \frac{(n - x + \beta - 1)}{1 - \theta} \\ &= \frac{x + \alpha - 1 - x\theta - \alpha\theta + \theta - n\theta + x\theta - \beta\theta + \theta}{\theta(1 - \theta)} \\ &= \frac{x + \alpha - 1 - (n + \alpha + \beta - 2)\theta}{\theta(1 - \theta)} = 0 \end{aligned}$$

$\theta(1 - \theta) \neq 0$ とすると

$$\hat{\theta} = \frac{x + \alpha - 1}{n + \alpha + \beta - 2}$$

がベイズ推定値となる。さて、 $\alpha, \beta$  は事前分布のパラメータであるが、これをハイパーパラメータ (hyper parameter) と呼ぶ。

ハイパーパラ  
メータによって、  
事前分布はさま  
ざまな形状をと  
る(図). 例えば、  
事前分布が一  
様となる場合  
( $Beta(1, 1)$ )の  
推定値は、 $\hat{\theta} =$   
 $\frac{x}{n}$ となり、最尤解  
に一致する.



# EAP推定量

$$\hat{\theta} = \frac{x + \alpha}{n + \alpha + \beta}$$

となり、例えば、事前分布が一様となる  
場合 ( $Beta(1, 1)$ ) の推定値は

$$\hat{\theta} = \frac{x + 1}{n + 2}$$

データがない場合は、 $\hat{\theta} = \frac{1}{2}$  となり、データが  
増えるごとに真値に近づく。

# EAP推定量でジェフリーズ事前分布

$$\hat{\theta} = \frac{x + \alpha}{n + \alpha + \beta}$$

となり、例えば、事前分布が一様となる  
場合 ( $Beta(1, 1)$ ) の推定値は

$$\hat{\theta} = \frac{x + 1/2}{n + 1}$$

データがない場合は、一様分布同様に  $\hat{\theta} = \frac{1}{2}$   
となるが、一様分布よりもデータに速く影響  
を受ける。

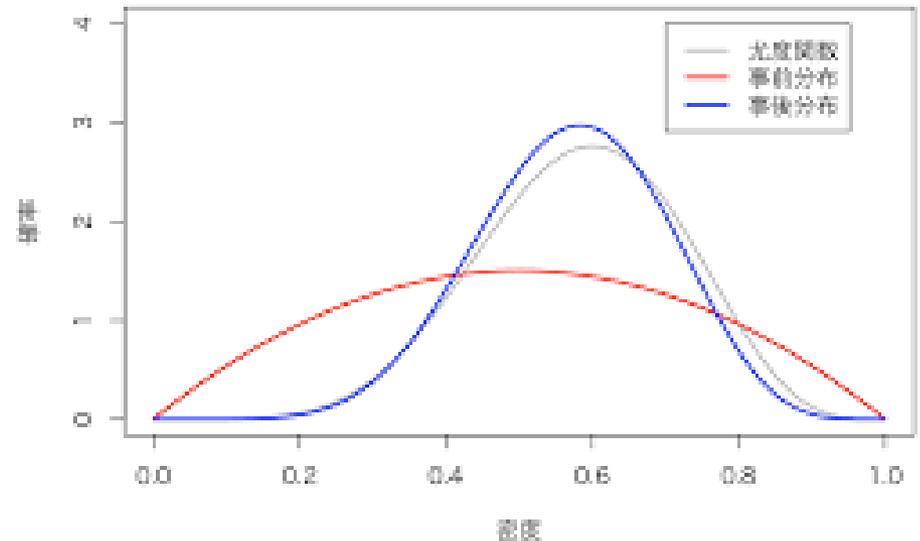
# 事後平均と事後分散

事後平均

$$\frac{\alpha + x}{\alpha + \beta + n}$$

事後分散

$$\frac{(\alpha + x)(\beta + n - x)}{(\alpha + \beta + n)^2(\alpha + \beta + n + 1)}$$



# 事後平均と事後分散

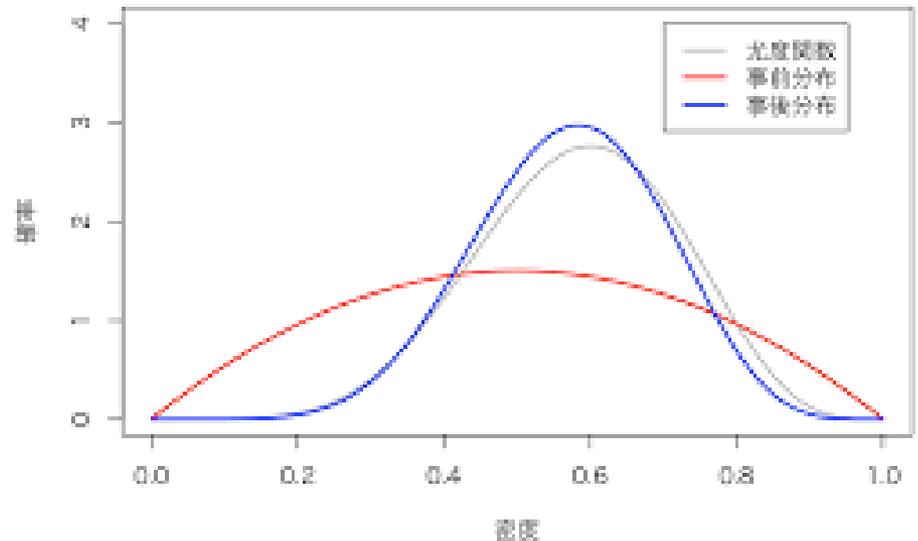
## 事後平均

$$\frac{\alpha + x}{\alpha + \beta + n}$$

## 事後分散

$$\frac{(\alpha + x)(\beta + n - x)}{(\alpha + \beta + n)^2(\alpha + \beta + n + 1)}$$

フィッシャー情報量の逆数はデータ数を無限を仮定。事後分散はデータが少なくてもその不確実性、あいまいさを反映。



## 例2 (正規分布)

$$P(x_i | \mu, \sigma^2) \\ = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\}$$

$(x_1, \dots, x_n)$  を得たときの  $\mu, \sigma^2$  を求めよう.

尤度は,

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$
$$= \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$

このとき, 自然共役事前分布は  $\sigma_0^2 = \frac{\sigma^2}{n_0}$  (注:  $n_0$  事前分布への信念の強さ)

$$p(\mu) = N(\mu_0, \sigma_0^2)$$
$$= \frac{1}{\sqrt{2\pi}\sigma_0} \exp \left\{ -\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right\}$$
$$\propto \left( \frac{\sigma^2}{n_0} \right)^{-\frac{1}{2}} \exp \left\{ -\frac{n_0(\mu - \mu_0)^2}{2\sigma^2} \right\}$$

$$p(\sigma^2) = Ig(\nu_0, \lambda_0)$$
$$= \frac{(\lambda_0/2)^{\frac{1}{2}\nu_0}}{\Gamma(\frac{1}{2}\nu_0)} (\sigma^2)^{-\frac{1}{2}\nu_0-1} \exp\left(-\frac{\lambda_0}{2\sigma^2}\right)$$

(逆ガンマ分布)

$$\propto (\sigma^2)^{-\frac{1}{2}\nu_0-1} \exp\left(-\frac{\lambda_0}{2\sigma^2}\right)$$

事前分布はこれらの積の形で以下のように表される. 自由度  $\nu_0 = n_0 - 1$  とすると

$$p(\mu, \sigma^2) = p(\mu | \mu_0, \sigma_0^2) p(\sigma^2 | \nu_0, \lambda_0)$$

$$\propto \left( \frac{\sigma^2}{n_0} \right)^{-\frac{1}{2}} \exp \left\{ -\frac{n_0(\mu - \mu_0)^2}{2\sigma^2} \right\}$$

$$(\sigma^2)^{-\frac{1}{2}\nu_0 - 1} \exp \left( -\frac{\lambda_0}{2\sigma^2} \right)$$

$$\propto (\sigma^2)^{-\frac{1}{2}(\nu_0 + 1) - 1} \exp \left\{ -\frac{\lambda_0 + n_0(\mu - \mu_0)^2}{2\sigma^2} \right\}$$

ここで  $n_0 = \nu_0 + 1$

事前分布を尤度に掛け合わせて事後分布を導くのであるが、計算の簡便さのために、以下のように尤度を変形させる。

$$L = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$

ここで指数部分  $\exp \left\{ - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right\}$  を三平方の定理により、推定平均  $\bar{x}$  を介して、以下のように分解する。

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{2\sigma^2} + \frac{(\bar{x} - \mu)^2}{2\sigma^2}$$

尤度 $L$ は、
$$L = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right\}$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{2\sigma^2}\right\}$$

$$\exp\left\{-\frac{n(\bar{x} - \mu)^2}{2\sigma^2}\right\}$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\frac{S^2 + n(\mu - \bar{x})^2}{2\sigma^2}\right\}$$

ただし、ここで、 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $S^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

$$\begin{aligned}
 p(\mu, \sigma^2 | x) &\propto L \times p(\mu, \sigma^2) \\
 &= \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{S^2 + n(\mu - \bar{x})^2}{2\sigma^2} \right\}
 \end{aligned}$$

$$\nu_0 = n_0 - 1 \text{ より}$$

$$\times (\sigma^2)^{-\frac{1}{2}(\nu_0+1)-1} \exp \left\{ -\frac{\lambda_0 + n_0(\mu - \mu_0)^2}{2\sigma^2} \right\}$$

$$\propto (\sigma^2)^{-\frac{1}{2}(n+n_0)-1}$$

$$\exp \left\{ -\frac{\lambda_0 + S^2 + n_0(\mu - \mu_0)^2 + n(\mu - \bar{x})^2}{2\sigma^2} \right\}$$

さらに、指数部分のうち、 $\lambda_0 + S^2$ 以外の部分に、平方完成を行うと、
$$p(\mu, \sigma^2 | x) \propto (\sigma^2)^{-\frac{1}{2}(n+n_0)-1} \exp \left\{ -\frac{\lambda_* + (n_0+n)(\mu-\mu_*)^2}{2\sigma^2} \right\}$$

ただし、
$$\lambda_* = \lambda_0 + S^2 + \frac{n_0 n (\bar{x} - \mu_0)^2}{n_0 + n}, \mu_* = \frac{n_0 \mu_0 + n \bar{x}}{n_0 + n}$$

この事後分布もまた、正規分布と逆ガンマ分布の積となり、

$$N \times IG(n_0 + n, \mu_*, \nu_0 + n, \lambda_*)$$

事後分布は、 $\mu$ と $\sigma^2$ の同時事後確率分布

# $\mu$ の周辺事後分布

このように、複数のパラメータを同時に最大化させる場合、つぎのような周辺化 (marginalization) を行い、個々のパラメータの分布を導く。このような分布を周辺事後分布 (marginal posterior distribution) と呼ぶ。  $p(\mu | x) = \int_0^\infty p(\mu, \sigma^2 | x) p(\sigma^2) d\sigma^2$

$$\propto \frac{\Gamma\left[\frac{(\nu_* + 1)}{2}\right]}{\sqrt{\frac{\nu_* \pi \lambda_*}{n_*}} \Gamma\left(\frac{\nu_*}{2}\right)} \left\{ 1 + \frac{(\mu - \mu_*)^2}{\mu_*} \right\}^{-\frac{1}{2}(\nu_* + 1)}$$
$$\equiv t(\nu_*, \mu_*, \lambda_*/n_*)$$

$\mu$  の周辺事後分布は  $t$  分布  $t(\nu_*, \mu_*, \lambda_*/n_*)$  に従う。

# MAP推定値

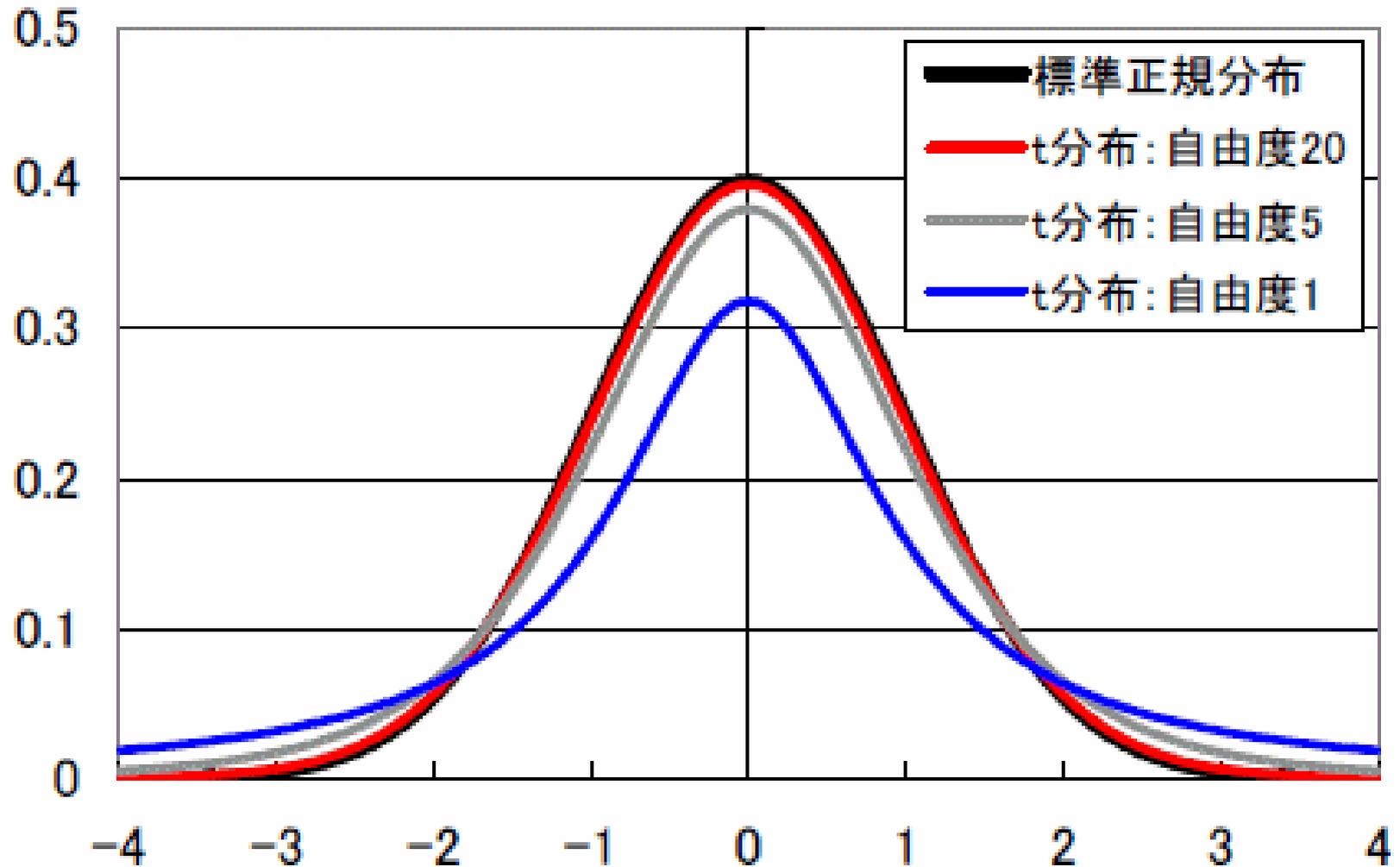
事後確率最大化によるベイズ推定値は、 $t$  分布のモードが  $\mu_*$  であることより、

$\mu$  のMAP推定値は、

$$\hat{\mu} = \frac{n_0 \mu_0 + n \bar{x}}{n_0 + n}$$

# 正規分布とt分布

## 標準正規分布とt分布



# $\sigma^2$ の周辺事後分布

$\sigma^2$ についての周辺事後分布は

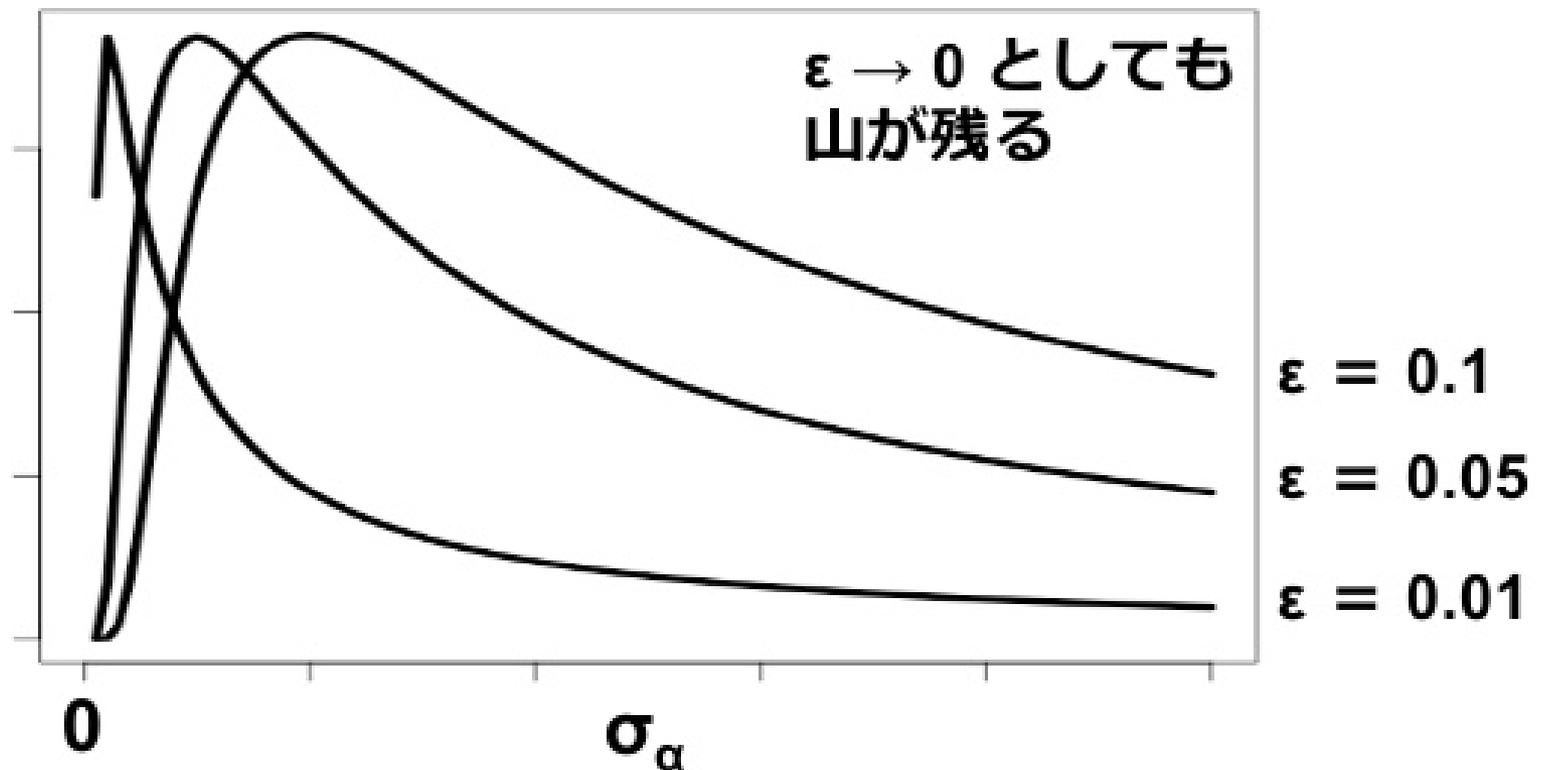
$$p(\sigma^2|x) = \int_0^{\infty} p(\mu, \sigma^2|x)p(\mu)d\mu$$

$$\propto \frac{\lambda_*^{\frac{\nu_*}{2}}}{2^{\frac{\nu_*}{2}}\Gamma\left(\frac{\nu_*}{2}\right)} (\sigma^2)^{-\frac{\nu_*}{2}-1} \exp\left(-\frac{\lambda_*}{2\sigma^2}\right)$$

となり,  $\sigma^2$ の周辺事後分布は, 逆ガンマ分布  $IG(\nu_*/2, \lambda_*/2)$  に従うことがわかる.

# 逆ガンマ分布

- $\sigma_\alpha \sim \text{InvGamma}(\varepsilon, \varepsilon)$



# MAP推定値

$\sigma^2$ のベイズ推定値は、逆ガンマ分布のモードが $\frac{\lambda_{*}/2}{\nu_{*}/2+1} = \frac{\lambda_{*}}{\nu_{*}+2}$ であることより、 $\sigma^2$ のMAP推定値は、

$$\widehat{\sigma^2} = \frac{\left\{ \lambda_0 + S^2 + \frac{n_0 n (\bar{x} - \mu_0)^2}{n_0 + n} \right\}}{\nu_* + 2}$$

# EAP推定値

$\mu$  のEAP推定値は，平均値とモードが同一なので

$$\hat{\mu} = \frac{n_0\mu_0 + n\bar{x}}{n_0 + n}$$

$\sigma^2$  のMAP推定値は，逆ガンマ分布のモードが  $\frac{\lambda_{*}/2}{\nu_{*}/2-1} = \frac{\lambda_{*}}{\nu_{*}-2}$  であることより，

$$\widehat{\sigma^2} = \frac{\left\{ \lambda_0 + S^2 + \frac{n_0 n (\bar{x} - \mu_0)^2}{n_0 + n} \right\}}{\nu_{*} - 2}$$

# 事後分散

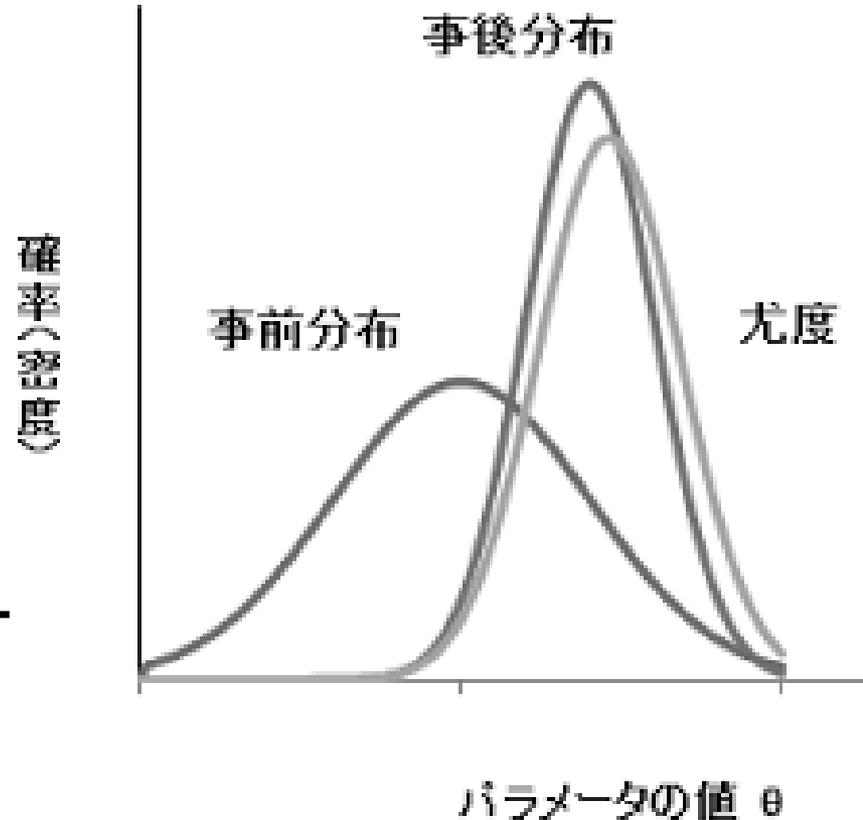
フィッシャー情報量と違い、データ数が少ない場合の不確実性を反映

平均値の事後分散

$$\frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

分散の事後分散

$$\frac{\lambda^*/2^2}{(v_{*/2} - 2)(v_{*/2} - 1)^2}$$



### 3. 事前分布の意味を考える例題

以下のどちらのかけを選ぶと得か？

1. 50個の赤玉と50個の白玉が入った壺から一つ玉を取り出し、それが赤玉であったら1万円もらえる。白玉であったら1万円支払う。
2. 赤玉と白玉が合わせて100個入った壺から一つ玉を取り出し、それが赤玉であったら1万円もらえる。白玉であったら1万円支払う。

1. の赤玉の出る確率は

1. 50個の赤玉と50個の白玉が入った壺から一つ玉を取り出し、それが赤玉(A)の確率

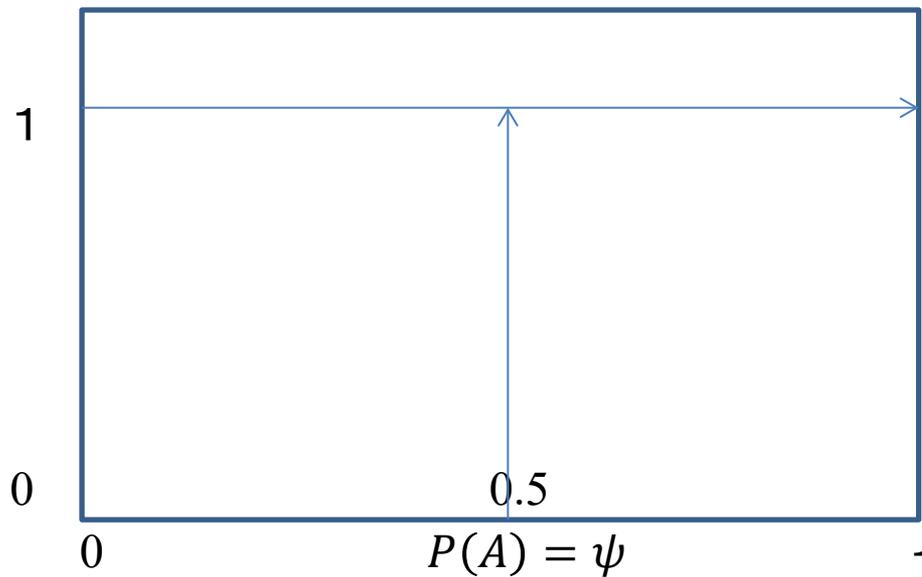
$$P(A) = \frac{50}{50 + 50}$$

## 2.の赤玉の出る確率は

2. 赤玉と白玉が合わせて100個入った壺から一つ玉を取り出し, それが赤玉の確率 $P(A) = \psi$ とする。

$$E(P(\psi)) = \int_0^1 \psi P(\psi) d\psi = \frac{1}{2}$$

確率 $P(A)$ の確率 $P(\psi)$



# 追加例題

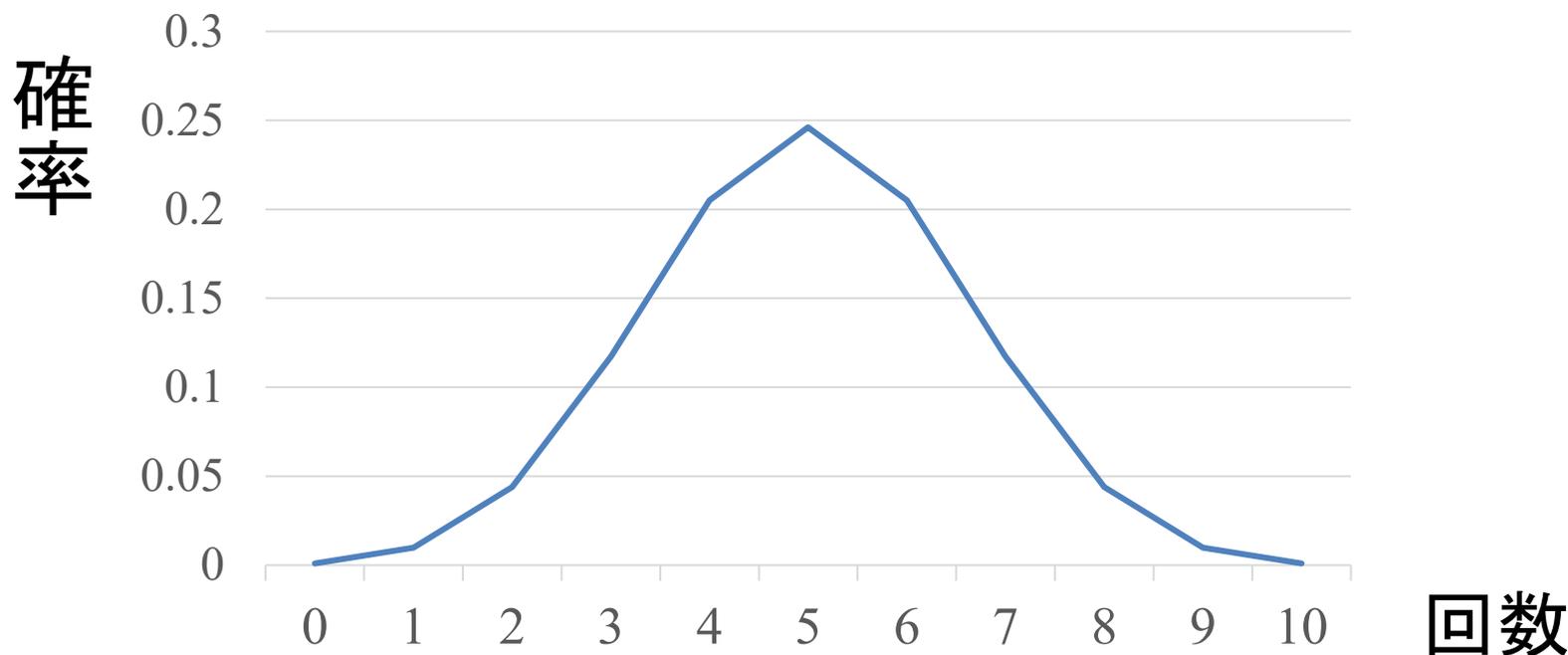
以下のどちらのかけを選ぶと得か？

1. 50個の赤玉と50個の白玉が入った壺から一つ玉を取り出し、それが赤玉であったら1万円もらえる。白玉であったら1万円支払う。これを10回繰り返す。
2. 赤玉と白玉が合わせて100個入った壺から一つ玉を取り出し、それが赤玉であったら1万円もらえる。白玉であったら1万円支払う。  
これを10回繰り返す。

# 分布を考えよう

1. 赤玉の出る回数を $x$ , 試行回数を $n$ としよう.  $p(x|\psi, n)$ は以下の二項分布に従う.

$$p(x|\psi, n) = \binom{n}{x} \psi^x (1 - \psi)^{n-x}$$

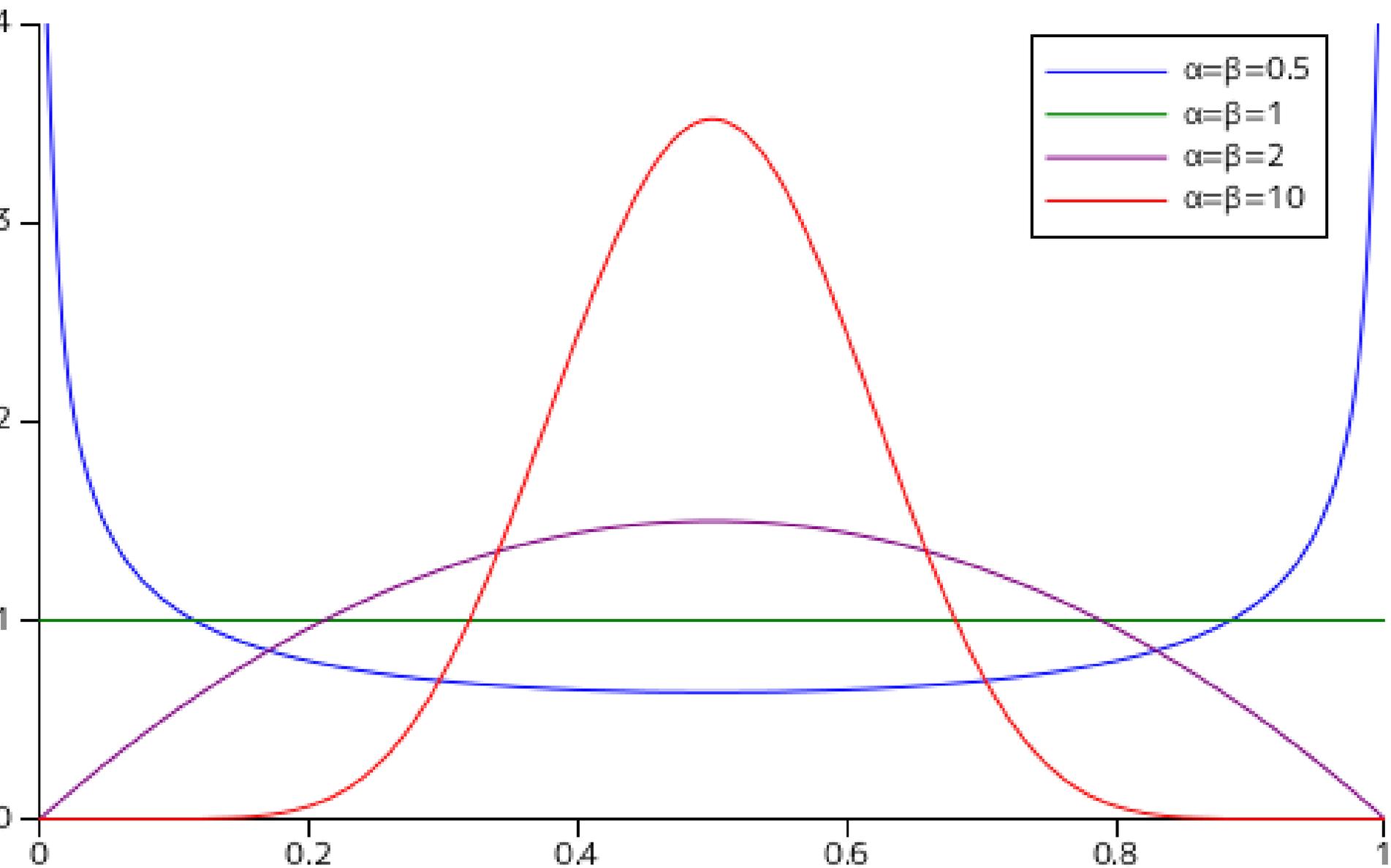


# 分布を考えよう

2. 赤玉の出る回数を $x$ , 試行回数を $n$ としよう. 事前分布をベータ分布とすると $p(x|\psi, n)$ は以下のベータ分布に従う.

$$p(\psi|n, x, \alpha, \beta) = \frac{\Gamma(n + \alpha + \beta)}{\Gamma(x + \alpha)\Gamma(n - x + \beta)} \psi^{x+\alpha-1} (1 - \psi)^{n-x+\beta-1}$$

問 ハイパーパラメータ $\alpha, \beta$  はどのように設定すればよいか?



赤玉の確率  $P(A) = \psi$

# 分布を考えよう

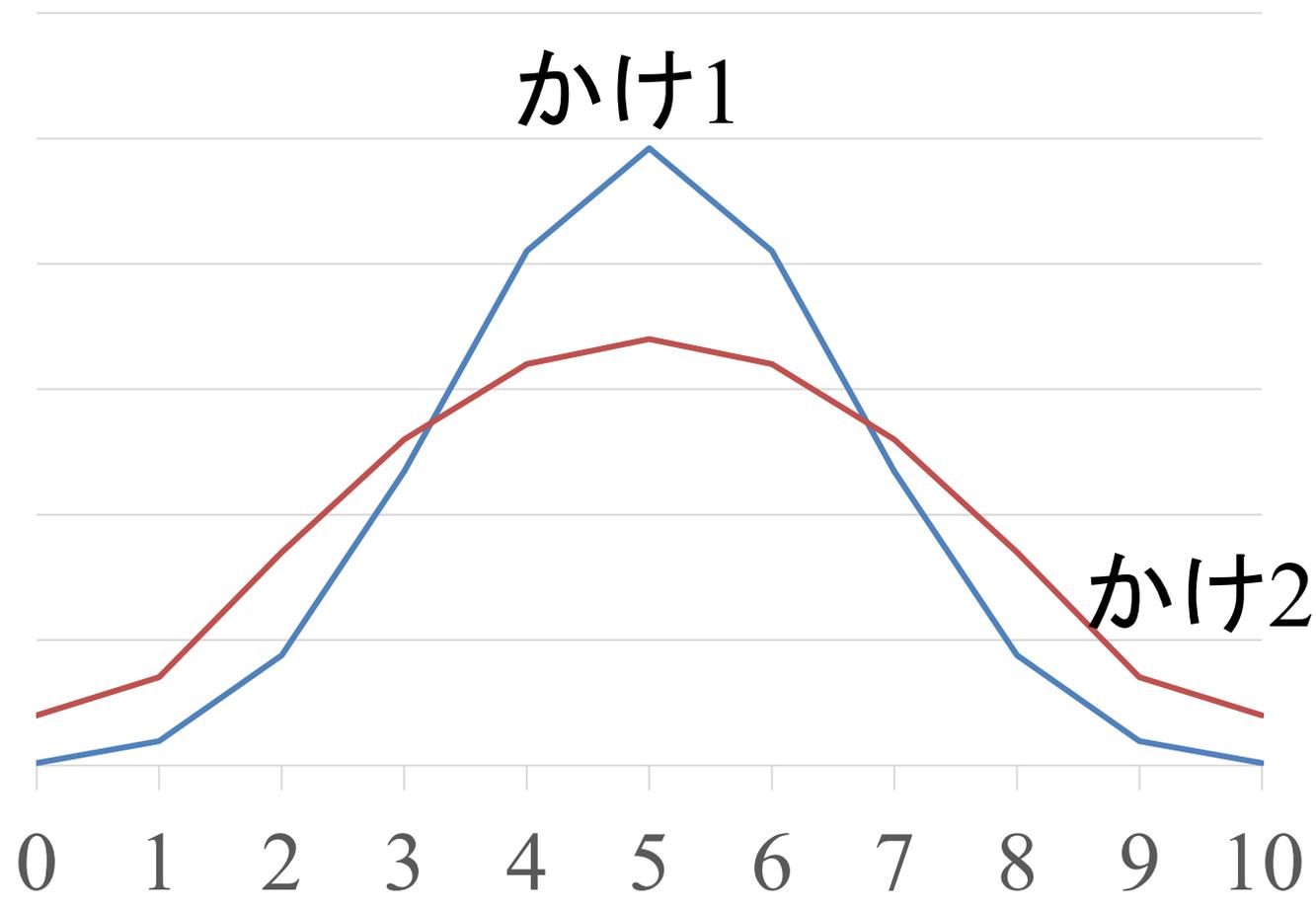
2. 赤玉の出る回数を $x$ , 試行回数を $n$ としよう. 事前分布をベータ分布とすると $p(x|\psi, n)$ は以下のベータ分布に従う.

$$\begin{aligned} p(\theta|n, x, \alpha = 1, \beta = 1) \\ &= \frac{\Gamma(n+2)}{\Gamma(x+1)\Gamma(n-x+1)} \theta^x (1-\theta)^{n-x} \end{aligned}$$

# 賭け1は博打性大

確率

0.3  
0.25  
0.2  
0.15  
0.1  
0.05  
0



かけ1

かけ2

回数

# 多くの人は「かけ1」を選ぶ

期待値が同じでも 多くの人は「かけ1」を選ぶことが知られている。

経済学でこの現象は人間の意思決定を予測する意味で重要である。

「人間は 利得よりも損失を過大評価する」ため 損失回避の方向で意思決定してしまうためであると解釈されている。

## 事前分布の例題2

いま、外見がまったく同じ2つの封筒の中に、現金が入っているものとする。それぞれの封筒の中の金額は知らされていないが、片方にはもう一方の2倍が入っていることが分かっている。今、AとBの二人に封筒がランダムに分けられ、自分の中身だけ見て交換してもよいルールとなった。Aの封筒には10ドル入っていた。交換したほうがよいか？

# 期待値を計算してみよう！！

自分は $X=10$ ドル入っていたので、相手は $Y=5$ ドルか $20$ ドルを持っている。その確率はそれぞれ $p(X)=p(Y)=0.5$ なので交換したときの期待値は  $5 \times 0.5 + 20 \times 0.5 = 10.25$ ドル。

今、持っているのは $10$ ドルなので交換したほうが良い！！

# 相手の立場になろう

相手はYドル持っていた場合もこちらが $X=1/2$  Yドルか  $X=2Y$ ドル 持っていることになる。同じ期待値の計算をすると  $0.5 \times 1/2 Y + 0.5 \times 2Y = 1.25Y$ ドルとなる。今 Yドル持っているので 交換したほうが得になる！！

え？

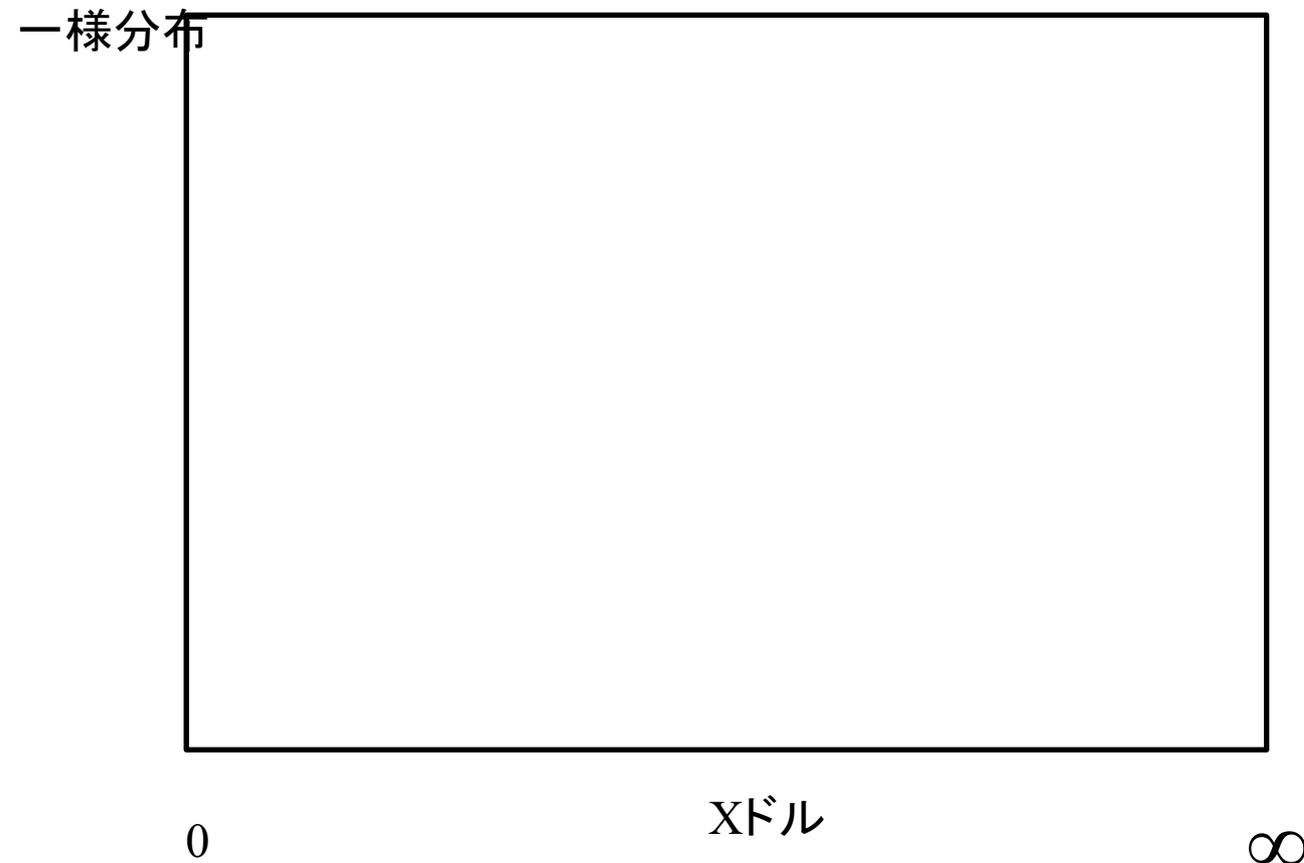
でも、相手も同じだよね。相手も 交換したほうが期待値が大きくなっているはず。。

どちらかが得すればどちらかが損するはずなのに、どちらも得するって変！！

なんで こんなことになるのでしょうか？

$$p(Y|X) = p(Y = 2X|X) = p(Y = \frac{1}{2}X|Y)$$

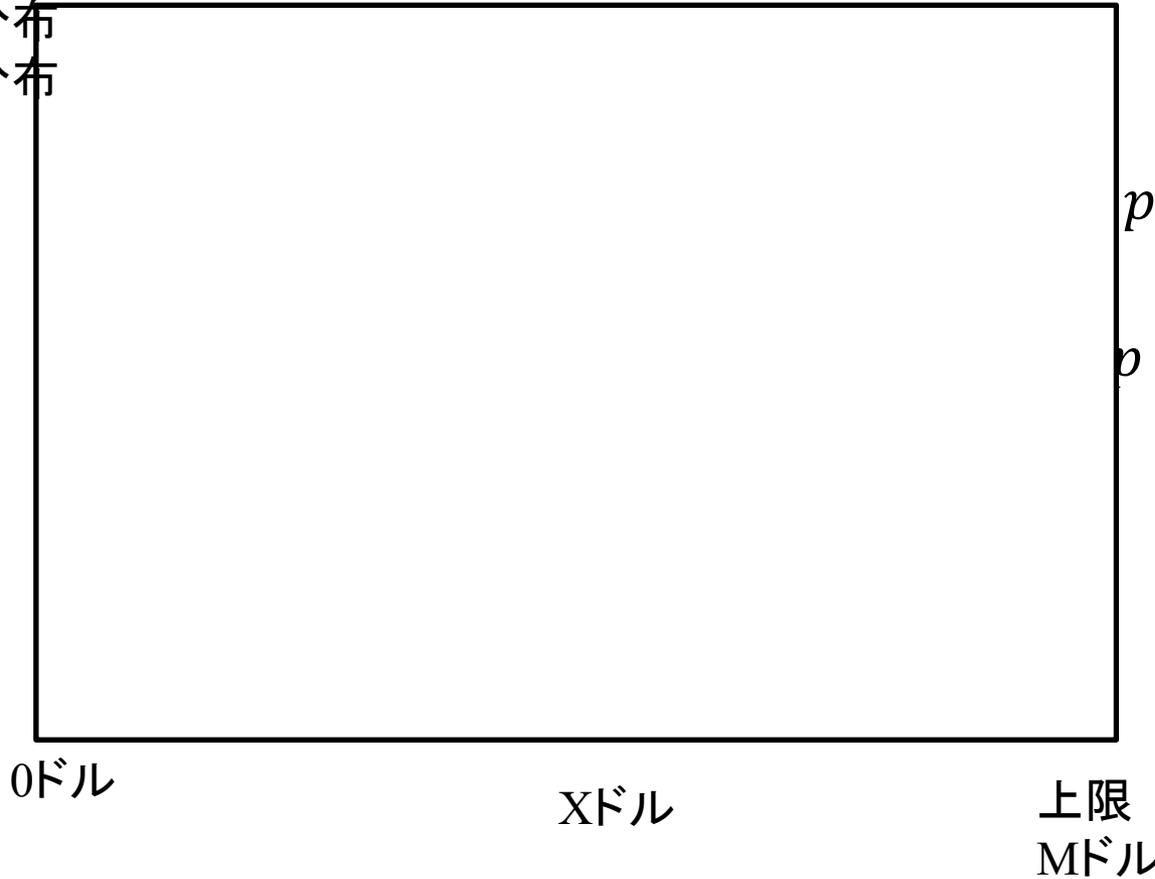
で暗黙に想定された事前分布



これは確率  
分布ではな  
い。

# 上限のある事前分布を考える

確率分布  
一様分布



$$p\left(Y = 2X \mid X \geq \frac{M}{2}\right) = 0$$
$$p\left(Y = 2X \mid X \leq \frac{M}{2}\right) = \frac{1}{2}$$

# 事前分布にガンマ分布を考える

$$p(X|\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} X^{\alpha-1} e^{-X/\beta}$$

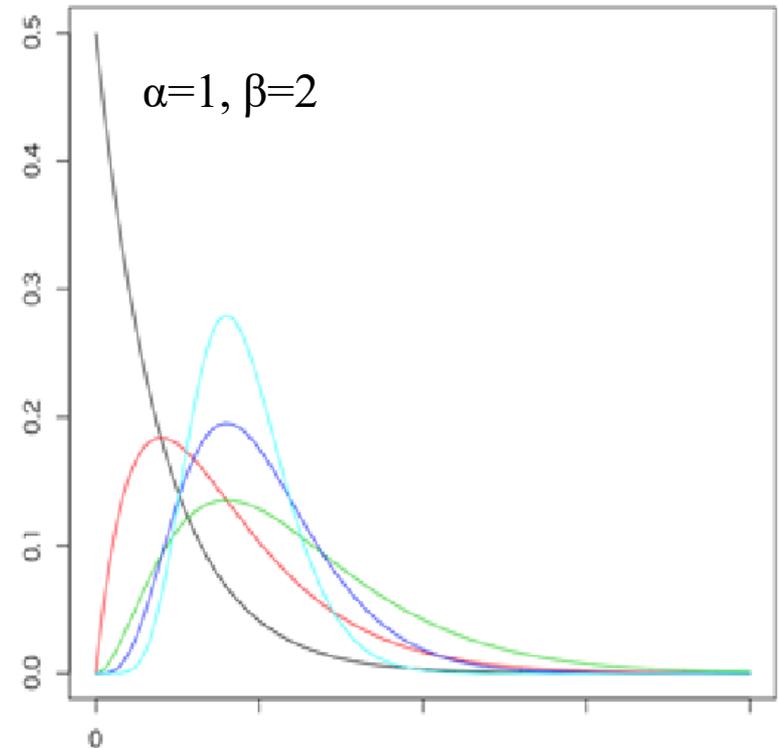
$$E(X) = \alpha\beta$$

$$X < 2(\alpha + 1)\beta \log 2 \\ \approx 0.6E(X)$$

のとき

$$E(Y|X) > X$$

交換すべき



## 4. データから統計モデルを選択

統計モデルのパラメータ(母数)をデータから推定するには、尤度最大化により漸近的な一致性が得られた。

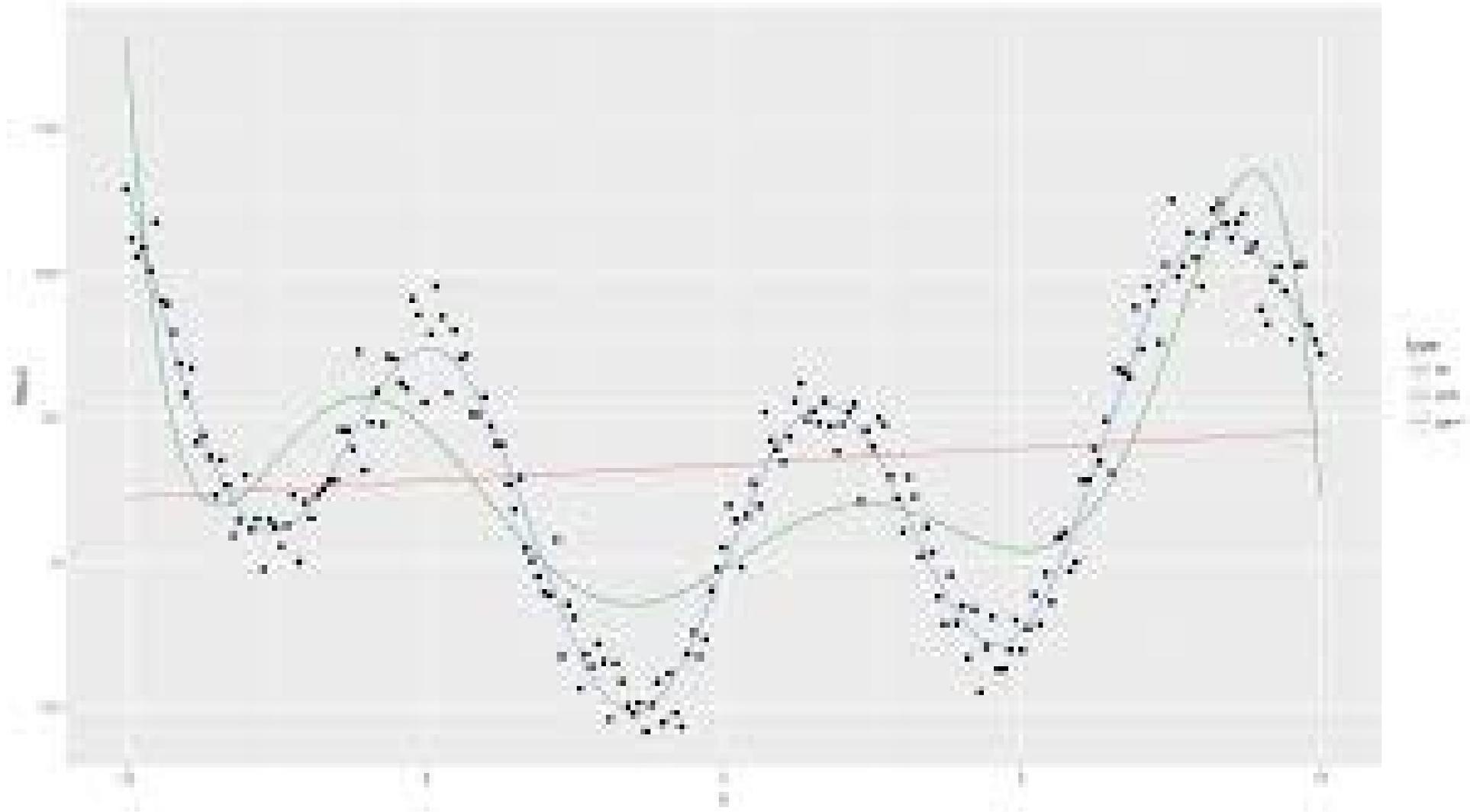
ひとつのデータに対して、複数のモデルからどのモデルが一番よいかを決定するとき、尤度最大化は使えるのでしょうか？

→

モデル選択基準

# 例；多項式のデータへのあてはめ

$$y = a_k x^k + a_{k-1} x^{k-1} + \dots * a_1 x + a_0$$



# パラメータ数が増えると予測が劣化

$$y = a_k x^k + a_{k-1} x^{k-1} + \dots * a_1 x + a_0$$

パラメータ数 =  $k+1$

パラメータ数が増える(モデルが複雑になる)とデータとの誤差が単調減少し、尤度は単調増加する。

データ数=パラメータ数のとき既知のデータへのあてはまり誤差は0になるが、未知のデータへの予測は非常に悪くなる。この現象を過学習(over fitting)という。

尤度最大化はモデル選択に使えない

複雑なモデルほど尤度が高くなってしまふので  
尤度最大化では、モデルの選択はできない

予測を最大にするモデルを選択手法は何か？

# ベイズではモデルの確率を考える

$m$ :モデル,  $M$ :モデル候補集合,  $x$ :データ

$$p(m|x) = \frac{p(x|m)p(m)}{\sum_{i=1}^M p(x|m_i)p(m_i)}$$

今、すべての $p(m)$ が同一だと考えると

$p(x|m)$  が最大となるモデルを選択すればよい。

ここで

$$p(x|m) = \int_{\Theta} p(x|\theta, m)p(\theta|m)d\theta$$

を周辺尤度と呼ぶ。

# 周辺尤度

ベイズ統計では、一般的に、モデル選択のために以下の周辺尤度を最大にするモデルを選択する。

定義19

データ $x$ を所与としたモデル $m$ の尤度を周辺化して周辺尤度 (marginal likelihood), **ML** と呼ぶ。

$$p(x|m) = \int_{\Theta} p(x|\theta, m)p(\theta|m)d\theta$$

# BIC(Bayesian Information Criterion)

周辺尤度は、モデルごとにパラメータ空間を積分消去しなければならない。より、簡単に用いるために 周辺尤度の漸近近似としてBICが求められた。これは漸近一致性を持つ。

$$\text{BIC} = \ln(L) - \frac{1}{2}k \ln(n)$$

ここで、 $\ln L$ は対数最大尤度、 $k$ はモデルのパラメータ数、 $n$ はデータ数。

Schwarz, Gideon E. (1978), "Estimating the dimension of a model", [\*Annals of Statistics\*](#), 6 (2): 461–464

# MDL(minimum description length)

Jorma Rissanen により導入された。MDLでは、データをモデルを用いて圧縮・送信する際の符号長の最小化を考える。これはノイズを含むデータから意味のある規則性を抽出することにあたる。最初はBICと等価な基準が提案されたが、その後NML( Normalized Maximum Likelihood )も提案されている。基本、符号問題、離散データの圧縮問題に用いる理論的仮定がある。

Rissanen, J. (1978). "Modeling by shortest data description". *Automatica*. **14** (5): 465–658

# MDL(minimum description length)

NMLのアイデアは尤度を確率になるように標準化する。そのためにはデータのとりえるパターンの尤度をすべて列挙して計算しなければならないので計算量の問題、またデータがスパースなパターンがあるのでデータスパース問題がありえる。

BICの数式は そもそも周辺尤度の近似であるが NMLの近似としても導ける。

# 世界最初の情報量基準AIC

## **Information of the NEW LOOK AT STATISTICAL-MODEL IDENTIFICATION**

By:AKAIKE, H (AKAIKE, H)

IEEE TRANSACTIONS ON AUTOMATIC  
CONTROL

Volume: AC19

Issue: 6, 1974

# AIC(Akaike Information Criterion 1973)

$$AIC = -2E[\ln L] \approx -2\ln L + 2k$$

ここで、 $\ln L$ は対数最大尤度、 $k$ はモデルのパラメータ数

Akaike, H., "Information theory and an extension of the maximum likelihood principle", *Proceedings of the 2nd International Symposium on Information Theory*, Petrov, B. N., and Caski, F. (eds.), Akadimiai Kiado, Budapest: 267-281 (1973).

# AICの意味

- $-\frac{1}{2} \text{AIC} =$  尤度 (モデルのてはまり)
- パラメータ数 (モデルの複雑さ)

モデルのてはまりと モデルの複雑さのトレードオフが存在する。

# AICは真のモデルに一致性を持たない

AICは真の事後分布、真の尤度に対して  
一致性を持つ。

汎化誤差最小化の意味ではベイズ的手法  
よりも優れている場合がある。

# 準備：分布 $P(\theta)$ と分布 $Q(\theta)$ の距離

カルバックライブラー距離

$$\int_{\theta} P(\theta) \log \frac{P(\theta)}{Q(\theta)} d\theta$$

# AICの導出の考え方

真の分布 $P^*(\theta)$ と分布の推定値 $P(\theta)$ のカルバツクライブラー距離

$$\begin{aligned} & \int_{\theta} P^*(\theta) \log \frac{P^*(\theta)}{P(\theta)} d\theta \\ &= \int_{\theta} P^*(\theta) \log P^*(\theta) d\theta \\ & \quad - \int_{\theta} P^*(\theta) \log P(\theta) d\theta \end{aligned}$$

# AICの導出の考え方

真の分布 $P^*(\theta)$ と分布の推定値 $P(\theta)$ のカルバック  
クライブラー距離

$$\int_{\theta} P^*(\theta) \log \frac{P^*(\theta)}{P(\theta)} d\theta$$

$$= \int_{\theta} P^*(\theta) \log P^*(\theta) d\theta$$

Const

$$- \int_{\theta} P^*(\theta) \log P(\theta) d\theta$$

ここだけ  
考えれば  
よい

クロス エントロピー

# AICの導出の考え方

$$\begin{aligned} & - \int_{\theta} P^*(\theta) \log P(\theta) d\theta \\ & \approx - \int_{\theta} P(\theta) \log P(\theta) d\theta \\ & \approx -E[\ln L] \end{aligned}$$

$\ln L$ を二回 テーラー近似し、

$$\begin{aligned} & -E[\ln L] \\ & \approx -\ln L + k \end{aligned}$$

これを最小化すればよい。

# 問題

$P^*(\theta)$ を $P(\theta)$ に置き換えてしまうとクロスエントロピーは厳密には真の分布との距離を反映していない。

# 赤池先生の偉大な貢献

AICは モデル選択が当てはまり(尤度)とパラメータ数(モデルの複雑さ)のトレードオフがあることを主張したことは モデル選択分野の創始者ともいえる。

その後、赤池先生は モデル選択の概念を持っていたBayes 手法に興味を持ち、ABIC の開発や日本におけるベイズ統計の発展に寄与された。

# 補足 汎化誤差

予測分布と真の分布とのカルバックライブラー距離を（ベイズ）汎化誤差と呼ぶ。真の分布に近い予測を行えるモデルを選択する。

しかし、真の分布が結局 わからないのでこの部分については 近似しないといけない。例えば、すべてのモデルを考えたモデル平均などを近似として選ぶ手法も考えられているが、計算量が大きすぎる問題もある。

AICは汎化誤差最小化モデルを選択する近似手法として解釈できる。

# 補足 AIC、BICと周辺尤度

Ueno, M. [Learning networks determined by the ratio of prior and data](#) The Twenty-Sixth Conference on Uncertainty in Artificial Intelligence (UAI), 598-605, 2010

AIC、BICも周辺尤度の一つとして解釈できる。

周辺尤度の事前分布を変化させると様々な情報量基準に変化する

例

1. パラメータの事前分布が一様分布 $\rightarrow$ BIC
2. パラメータの事前分布が真の分布に一致 $\rightarrow$ AIC

# 5. 予測分布

データやモデルを用いて推論を行う重要な目的の一つに、未知の事象の予測が挙げられる。この予測問題のためには、最もよく用いられるのは、

$$p(y|\hat{\theta})$$

で示されるplug-in distribution と呼ばれる分布である。しかし、 $\hat{\theta}$  は推定値であるためにそのサンプルのとり方によってこの分布は大きく変化する。ベイズ的アプローチでは、この $\hat{\theta}$ のばらつき( $\hat{\theta}$ の事後分布)を考慮し、以下のように予測分布を定義する。

# 予測分布

## 定義3

モデル $m$  から発生されるデータ $x$  により, 未知の変数 $y$  の分布を予測するとき, 以下の分布を予測分布 (predictive distribution) と呼ぶ.

$$p(y|x, m) = \int_{\Theta} p(y|\theta, m)p(\theta|x, m)d\theta$$

例3 (二項分布) ベータ分布を事前分布とした二項分布の予測分布は, 以下のようになる.

$$\begin{aligned} p(y|x) &= \int_{\Theta} p(y|\theta)p(\theta|x)d\theta = \\ &\int_{\Theta} \binom{n}{y} \theta^y (1-\theta)^{n-y} \frac{\Gamma(n+\alpha+\beta)}{\Gamma(x+\alpha)\Gamma(n-x+\beta)} \times \\ &\theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta \\ &\propto \binom{n}{y} \frac{\Gamma(y+1)\Gamma(n-y+1)}{\Gamma(n+2)} \\ &\quad \frac{\Gamma(x+\alpha)(n-x+\beta)}{\Gamma(n+\alpha+\beta)} \end{aligned}$$

$$= \frac{n!}{y! (n-y)!} \frac{\Gamma(y+1)\Gamma(n-y+1)}{\Gamma(n+2)} \frac{\Gamma(x+\alpha)(n-x+\beta)}{\Gamma(n+\alpha+\beta)}$$

特に， $\alpha, \beta$  が整数のとき

$$p(y|x) \propto \frac{n!}{y! (n-y)!} \frac{y! (n-y)!}{(n+1)!} \frac{(x+\alpha-1)! (n-x+\beta-1)!}{(n+\alpha+\beta-1)!}$$

例4 (正規分布) 事前分布を  $N(\mu, \sigma^2)$  分布

$$p(\mu, \sigma^2) = p(\mu|\sigma^2)p(\sigma^2)$$

$$\begin{aligned} &\propto \left(\frac{\sigma^2}{n_0}\right)^{-\frac{1}{2}} \exp\left\{-\frac{n_0(\mu - \mu_0)^2}{2\sigma^2}\right\} (\sigma^2)^{-\frac{1}{2}\nu_0-1} \\ &\quad \exp\left(-\frac{\lambda_0}{2\sigma^2}\right) \\ &= (\sigma^2)^{-\frac{1}{2}(\nu_0+1)-1} \exp\left\{-\frac{\lambda_0 + n_0(\mu - \mu_0)^2}{2\sigma^2}\right\} \end{aligned}$$

事後分布は

$$p(\mu, \sigma^2 | x) \propto (\sigma^2)^{-\frac{1}{2}(n+n_0)-1} \exp \left\{ -\frac{\lambda_* + (n_0 + n)(\mu - \mu_*)^2}{2\sigma^2} \right\}$$

ただし,  $\lambda_* = \lambda_0 + S^2 + \frac{n_0 n (\bar{x} - \mu_0)^2}{n_0 + n},$

$$\mu_* = \frac{n_0 \mu_0 + n \bar{x}}{n_0 + n}$$

予測分布は

$$p(x_{n+1} | \mathbf{x})$$

$$= \int \int p(x_{n+1} | \mu, \sigma^2) p(\mu, \sigma^2 | x_1, \dots, x_n) d\mu d\sigma^2$$

ここで,  $p(x_{n+1} | \mu, \sigma^2) \propto$   
 $(\sigma^2)^{-1} \exp \left\{ -\frac{(x_{n+1} - \mu)^2}{2\sigma^2} \right\}$

$$\begin{aligned}
p(x_{n+1}|\mathbf{x}) &= \int \int p(x_{n+1}|\mu, \sigma^2)p(\mu, \sigma^2|x_1, \dots, x_n)d\mu d\sigma^2 \\
&\propto \int \int (\sigma^2)^{-\frac{\nu+1}{2}-2} \exp \left[ -\frac{(x_{n+1}-\mu)^2 + S^2 + n(\mu - \bar{x})^2}{2\sigma^2} \right] d\mu d\sigma^2 \\
&= \int \int (\sigma^2)^{-\frac{\nu+1}{2}-2} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (n+1)(\mu - \bar{\mu})^2 + S^2 \right. \right. \\
&\quad \left. \left. + \frac{n}{n+1}(x_{n+1} - \bar{x})^2 \right\} \right] d\mu d\sigma^2 \\
&\propto \int (\sigma^2)^{-\frac{\nu+1}{2}-2} \exp \left[ -\frac{1}{2\sigma^2} \left\{ S^2 + \frac{n}{n+1}(x_{n+1} - \bar{x})^2 \right\} \right] d\sigma^2 \\
&\propto \left\{ S^2 + \frac{n}{n+1}(x_{n+1} - \bar{x})^2 \right\}^{-\frac{\nu+1}{2}}
\end{aligned}$$

$$\propto \left[ 1 + \left\{ \frac{x_{n+1} - \bar{x}}{\sqrt{\frac{n+1}{nv} S^2}} \right\}^2 / \nu \right]^{-\frac{\nu+1}{2}}$$

ただし,ここで

$$\bar{\mu} = \frac{n\bar{x} + x_{n+1}}{n+1}$$

ここで,  $t = \frac{x_{n+1} - \bar{x}}{\sqrt{\frac{n+1}{nv} S^2}}$

とおくとき,  $t$ は自由度 $\nu$ の $t$ 分布に従う.

# 5. マルコフ連鎖モンテカルロ法 (MCMC法)

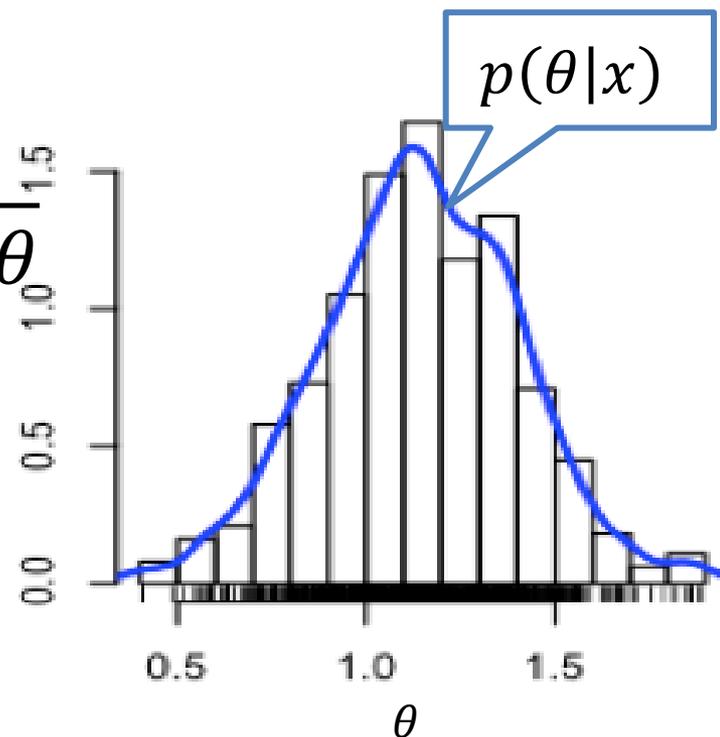
確率分布をサンプリング近似する手法

ベイズ推定では、パラメータの事後分布を推定し、  
得られた分布形に基づいて推定値を求める

$$p(\theta|\mathbf{x}) = \frac{p(\theta) \prod_{i=1}^n f(x_i|\theta)}{\int_{\Theta} p(\theta) \prod_{i=1}^n f(x_i|\theta) d\theta}$$

$\operatorname{argmax}_{\theta} p(\theta|\mathbf{x}) \rightarrow \text{MAP推定値}$

$E_{\theta} [ p(\theta|\mathbf{x}) ] \rightarrow \text{EAP推定値}$



# 代表的なMCMCアルゴリズム

1. ギブスサンプリング
2. メトロポリスヘイスティングス

他のMCMCアルゴリズム:

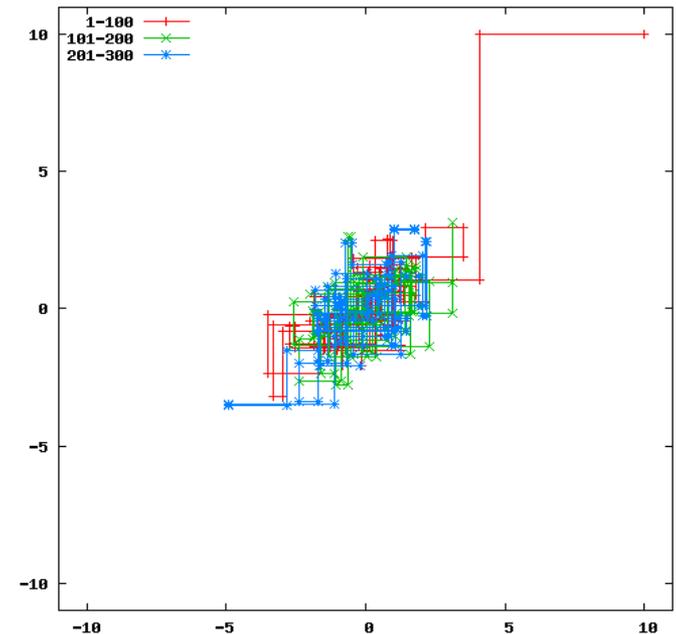
スライスサンプリング

ハミルトニアンモンテカルロ

以降では, 多次元パラメータ $\theta = \{\theta_1, \dots, \theta_K\}$ の事後分布をMCMCで推定することを想定

# 5.1. ギブスサンプリング

事後分布  $p(\theta|x)$  から直接にはサンプリングできないが、パラメータごとの条件付き分布  $p(\theta_i|x, \theta \setminus i)$  からサンプリングができる場合に利用できる手法(ここで、 $\theta \setminus i = \theta \setminus \{\theta_i\}$ )パラメータごとの条件付き分布から順にサンプリングを繰り返す



2次元正規分布の例

<http://d.hatena.ne.jp/jetbead/20120119/1326987540> より

# アルゴリズム

以下を十分な回数繰り返す

$$\theta_1 \sim p(\theta_1 | x, \boldsymbol{\theta} \setminus^1)$$

$$\theta_2 \sim p(\theta_2 | x, \boldsymbol{\theta} \setminus^2)$$

⋮

$$\theta_K \sim p(\theta_K | x, \boldsymbol{\theta} \setminus^K)$$

サンプリングしたパラメータ値  $\theta$  を保存

# 例：正規分布のパラメータ推定

$x_i \sim N(\mu, \sigma^2)$ とする $n$ 個のサンプル $\mathbf{x} = \{x_1, \dots, x_n\}$ を所与としてパラメータ $\mu, \sigma^2$ を推定  
パラメータの同時事後分布はサンプリング可能な既知の分布とならないため、この分布から直接サンプリングすることはできない

$$p(\mu, \sigma^2 | \mathbf{x}) = \frac{p(\mu)p(\sigma^2) \prod_{i=1}^n f(x_i | \mu, \sigma)}{\int p(\mu)p(\sigma^2) \prod_{i=1}^n f(x_i | \mu, \sigma) d\mu, \sigma}$$

しかし、以下の条件付き分布はそれぞれ既知の分布になるため、サンプリングが可能

$$p(\mu | \mathbf{x}, \sigma^2), p(\sigma^2 | \mathbf{x}, \mu)$$

$\mu, \sigma^2$ の事前分布に一様分布を仮定すると

$$p(\mu|\mathbf{x}, \sigma^2) = N\left(\frac{1}{N} \sum_{i=1}^n x_i, \frac{\sigma^2}{N}\right)$$

$$p(\sigma^2|\mathbf{x}, \mu) = IG\left(\frac{n}{2} + 1, \frac{\sum_{i=1}^n (x_i - \mu)^2}{2}\right)$$

正規分布や逆ガンマ分布 $IG()$ からの乱数生成手法は既知

多くのプログラミング言語にはこれらの乱数生成器が実装されている

## 5.2. メトロポリスヘイスティングス 条件付き分布からもサンプリングできないと きに利用

Step1:

現在のパラメータ値 $\theta$ の付近の候補値 $\theta^*$ を,  
提案分布 (proposal distribution)  $p(\theta^*|\theta)$  から生成

# 一般に  $q(\theta^*|\theta) = MN(\theta^*|\theta, I\sigma)$

MNは多次元正規分布,  $I$ は単位行列,  $\sigma$   
は微小な値(0.01等)

## 5.2. メトロポリスヘイスティングス

Step2:

以下の採択確率に基づいて候補値 $\theta^*$ を採択

$$\alpha(\theta^*, \theta) = \min \left\{ 1, \frac{p(\theta^* | \mathbf{x}) q(\theta | \theta^*)}{p(\theta | \mathbf{x}) q(\theta^* | \theta)} \right\}$$

( $q(\theta^* | \theta)$

=  $MN(\theta^* | \theta, I\sigma)$  のとき)  $\alpha(\theta^*, \theta)$

$$= \min \left\{ 1, \frac{p(\theta^* | \mathbf{x})}{p(\theta | \mathbf{x})} \right\}$$

棄却された場合には $\theta^* = \theta$ とする

# 考え方

$$p(\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*) = q(\boldsymbol{\theta}^* | \boldsymbol{\theta}) \alpha(\boldsymbol{\theta}^*)$$

$$p(\boldsymbol{\theta}^* \rightarrow \boldsymbol{\theta}) = q(\boldsymbol{\theta} | \boldsymbol{\theta}^*) \alpha(\boldsymbol{\theta})$$

$$\frac{\alpha(\boldsymbol{\theta}^*)}{\alpha(\boldsymbol{\theta})} = \frac{p(\boldsymbol{\theta}^* | \boldsymbol{x}) q(\boldsymbol{\theta} | \boldsymbol{\theta}^*)}{p(\boldsymbol{\theta} | \boldsymbol{x}) q(\boldsymbol{\theta}^* | \boldsymbol{\theta})}$$

# 採択確率計算時のポイント

事後分布の分母は多重積分を含むため計算困難

$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\boldsymbol{\theta}) \prod_{i=1}^n f(x_i|\boldsymbol{\theta})}{\int p(\boldsymbol{\theta}) \prod_{i=1}^n f(x_i|\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

しかし、採択確率の計算ではこの項は消去可能

$$\begin{aligned} \frac{p(\boldsymbol{\theta}^*|\mathbf{x})}{p(\boldsymbol{\theta}|\mathbf{x})} &= \frac{\frac{p(\boldsymbol{\theta}^*) \prod_{i=1}^n f(x_i|\boldsymbol{\theta}^*)}{\int p(\boldsymbol{\theta}^*) \prod_{i=1}^n f(x_i|\boldsymbol{\theta}^*) d\boldsymbol{\theta}^*}}{\frac{p(\boldsymbol{\theta}) \prod_{i=1}^n f(x_i|\boldsymbol{\theta})}{\int p(\boldsymbol{\theta}) \prod_{i=1}^n f(x_i|\boldsymbol{\theta}) d\boldsymbol{\theta}}} \\ &= \frac{p(\boldsymbol{\theta}^*) \prod_{i=1}^n f(x_i|\boldsymbol{\theta}^*)}{p(\boldsymbol{\theta}) \prod_{i=1}^n f(x_i|\boldsymbol{\theta})} \end{aligned}$$

## 5.3. メトロポリス with ギブス

メトロポリスヘイスティングスでは、パラメータ数が増加すると、パラメータ値が改悪される方向に進むときに、採択確率  $\frac{p(\boldsymbol{\theta}^*|\mathbf{x})}{p(\boldsymbol{\theta}|\mathbf{x})}$  が極端に小さくなり、更新が進まなくなることがある

### メトロポリスヘイスティングス with ギブス

ギブスサンプリングとのハイブリッド法. パラメータごとに他のパラメータの条件付分布を求めてメトロポリスヘイスティングスを実行する手法

# アルゴリズム

Init  $\boldsymbol{\theta} = \{\theta_1 \cdots \theta_K\}$  # 初期値をランダムに設定

For loop = 1 to Max Loop:

For  $i = 1, \dots, K$ :

- ・現在の値を所与として  $\theta_i$  の候補値  $\theta_i^*$  を生成

$$\theta_i^* \sim N(\theta_i, \sigma^2)$$

- ・採択確率に基づき  $\theta_i^*$  を採択 (または棄却)

$$\alpha(\theta_i^*, \theta_i) = \min \left\{ 1, \frac{p(\theta_i^* | x, \boldsymbol{\theta} \setminus i)}{p(\theta_i | x, \boldsymbol{\theta} \setminus i)} \right\}$$

end for

現在のパラメータ値を保存

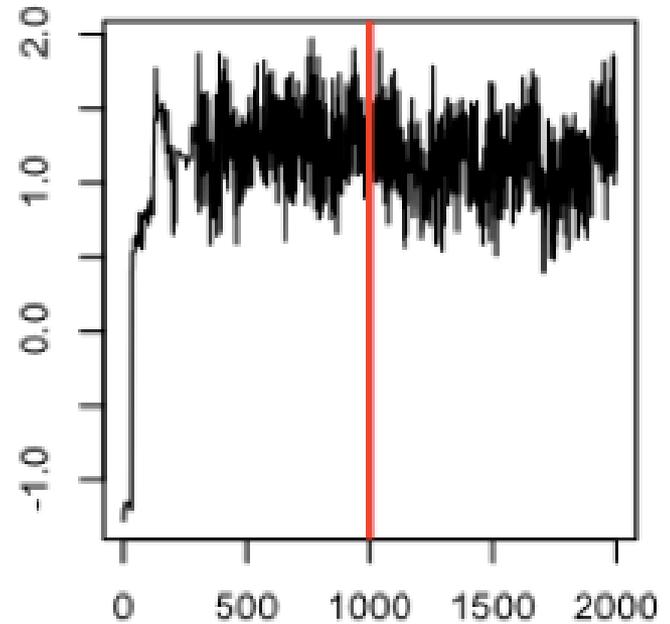
end for

得られたサンプル集合を事後分布からのサンプルとみなし、  
所望の統計量を計算

# サンプルの選択(バーンイン)

アルゴリズム初期のサンプルは、初期値に依存するため、一定回数サンプリングを繰り返した後のサンプルを利用する

初期値に依存しなくなったとみなすまでの時間をバーンイン期間と呼ぶ



# サンプルの選択（インターバル）

メトロポリスヘイスティングスは、サンプル間の自己相関（隣接するサンプル間の依存性）が高いため、一定区間でサンプルを間引いて用いる必要がある

間引く間隔をインターバル期間と呼ぶ

# アルゴリズム (修正版)

Init  $\theta = \{\theta_1 \cdots \theta_K\}$  # 初期値をランダムに設定

For loop = 1 to Max Loop:

For  $i = 1, \dots, K$ :

- ・現在の値を所与として  $\theta_i$  の候補値  $\theta_i^*$  を生成

$$\theta_i^* \sim N(\theta_i, \sigma^2)$$

- ・採択確率に基づき  $\theta_i^*$  を採択 (または棄却)

$$\alpha(\theta_i^*, \theta_i) = \min \left\{ 1, \frac{p(\theta_i^* | x, \theta \setminus i)}{p(\theta_i | x, \theta \setminus i)} \right\}$$

end for

If loop > **バーンイン期間**:

If loop % interval = 0:

**現在のパラメータ値  $\theta$  を保持**

end for

得られたサンプル集合を用いて所望の統計量を計算

# 6. 変分ベイズ学習の枠組み

ベイズ学習を汎関数の最小化問題として定式化

– 汎関数(functional): 関数を変数として持つ関数

変分自由エネルギー(variational free energy)

$$F(q) = \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)}{p(X|\theta)p(\theta)} \right]$$

$q(\theta)$ (または $q$ と略す): パラメータ $\theta$ の空間の任意の確率分布

$$F(q) = \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)}{p(X|\theta)p(\theta)} \right]$$

$$= \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)}{p(\theta|X)} \right] - \log p(X) \quad \leftarrow \text{ベイズ自由エネルギー}$$

$$= KL(q(\theta) || p(\theta|X)) + F^*$$

Kullback-Leibler divergence: 二つの分布の類似距離

$$KL(p_1(\theta) || p_2(\theta)) = \mathbb{E}_{p_1(\theta)} \left[ \log \frac{p_1(\theta)}{p_2(\theta)} \right] = \int p_1(\theta) \log \frac{p_1(\theta)}{p_2(\theta)} d\theta$$

自由エネルギー  $F(q)$  の最小化



$KL(q(\theta) || p(\theta|X))$  の最小化

$$F(q) = KL(q(\theta)||p(\theta|X)) + F^* \\ \geq F^*$$

⇒自由エネルギー $F(q)$ がベイズ自由エネルギー $F^*$   
 $= -\log p(X)$ の上界

自由エネルギーの符号反転:  $-F(q) \leq \log p(X)$

### 証拠の下界(evidence lower bound(ELBO))

自由エネルギー $F(q)$ の最小化  
ELBOの最大化



$KL(q(\theta)||p(\theta|X))$ の最小化

# 変分ベイズ学習の定式化

制約なし最小化問題

$$\hat{q} = \operatorname{argmin}_q F(q)$$

$\operatorname{argmin}_q F(q)$ の解はベイズ事後分布に一致

$$\hat{q} = p(\theta|X)$$

独立性の制約により自由エネルギーの期待値計算が可能に

**変分ベイズ事後分布**(variational Bayesian posterior):

$$\hat{q} = \underset{q}{\operatorname{argmin}} F(q), \quad \text{s.t. } q(\theta) = \prod_{s=1}^S q_s(\theta_s)$$

独立性の制約によって事後分布の各因子 $\{q_s\}_{s=1}^S$ を別々に最適化することが可能

⇒ 各因子を**変分法**を用いて最適化

# 変分法(calculus of variations)

汎関数の極値条件から解である関数が満たすべき条件を求める

**変分:** 変数関数 $q$ の微小変化に対する目的汎関数 $F(q)$ の変化量

$F(q)$ を最小とするような $q$

$\Rightarrow$ すべての $\theta$ の取りうる値に対して変分 $\delta I$ が0

$$\delta I = \frac{\partial F}{\partial q} = 0, \quad \forall \theta \in \Theta \quad (6.11)$$

# 変分ベイズ学習アルゴリズムの導出

実際に変分法を自由エネルギー最小化問題に  
適用する

自由エネルギー $F(r)$ に $q(\theta) = \prod_{s=1}^S q_s(\theta_s)$ ,  $p(\theta) = \prod_{s=1}^S p(\theta_s)$ を代入

$$\begin{aligned} F(q) &= \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)}{p(X|\theta)p(\theta)} \right] \\ &= \int q(\theta) \log \frac{q(\theta)}{p(X|\theta)p(\theta)} d\theta \\ &= \int \left( \prod_{s=1}^S q_s(\theta_s) \right) \left( \log \frac{\prod_{s=1}^S q_s(\theta_s)}{p(X|\theta) \prod_{s=1}^S p(\theta_s)} \right) d\theta \end{aligned}$$

$$\prod_{s=1}^S q_s(\theta_s) = q_s(\theta_s) \prod_{s' \neq s} q_{s'}(\theta_{s'})$$

と分解し各因子 $q_s(\theta_s)$ に対する自由エネルギーの変分を計算

$$\begin{aligned} 0 = \frac{\partial F}{\partial q_s} &= \int \left( \prod_{s' \neq s} q_{s'}(\theta_{s'}) \right) \left( \log \frac{\prod_{s'=1}^S q_{s'}(\theta_{s'})}{p(X|\theta) \prod_{s'=1}^S p(\theta_{s'})} + 1 \right) d\theta \\ &= \mathbb{E}_{\prod_{s' \neq s} q_{s'}(\theta_{s'})} \left[ \log \frac{\prod_{s'=1}^S q_{s'}(\theta_{s'})}{p(X|\theta) \prod_{s'=1}^S p(\theta_{s'})} \right] + 1 \\ &= \mathbb{E}_{\prod_{s' \neq s} q_{s'}(\theta_{s'})} \left[ \log \frac{\prod_{s' \neq s} q_{s'}(\theta_{s'})}{p(X|\theta) \prod_{s' \neq s} p(\theta_{s'})} \right] + \log \frac{q_s(\theta_s)}{p(\theta_s)} + 1 \\ &= \mathbb{E}_{\prod_{s' \neq s} q_{s'}(\theta_{s'})} \left[ \log \frac{1}{p(X|\theta)} \right] + \log \frac{q_s(\theta_s)}{p(\theta_s)} + \text{const} \end{aligned}$$

# 局所変分ベイズ学習アルゴリズム

$$\begin{aligned}\frac{\partial F}{\partial q_s} &= \mathbb{E}_{\prod_{s' \neq s} q_{s'}(\theta_{s'})} \left[ \log \frac{1}{p(X|\theta)} \right] + \log \frac{q_s(\theta_s)}{p(\theta_s)} + \text{const} \\ &= 0\end{aligned}$$

$$q_s(\theta_s) \propto p(\theta_s) \exp \left( \mathbb{E}_{\prod_{s' \neq s} q_{s'}(\theta_{s'})} [\log p(X|\theta)] \right)$$

# 変分ベイズ学習の各種計算量

変分ベイズ事後分布の平均値

$$\hat{\theta} = \mathbb{E}_{\hat{q}(\theta)}[\theta]$$

は**変分ベイズ推定量**(variational Bayesian estimator)と呼ばれる

変分ベイズ事後分布は正規分布などの関数形になるの  
で、事後平均および事後分散の計算は容易

# 7. (従来手法) 統計的仮説検定

- ある仮説が正しいかどうかを標本（データ）から判定する手法.
- 統計的仮説 (Statistical Hypothesis) :
  - 帰無仮説 (null hypothesis) : 棄却されることを前提とした仮説を表し  $H_0$  とする.
  - 対立仮説 (alternative hypothesis) : 帰無仮説が棄却されたときの採用される仮説を表し  $H_1$  とする.
- 有意水準  $\alpha$  : ユーザが設定する帰無仮説を棄却する基準であり, 誤って帰無仮説を棄却してしまう確率を表す.

# 仮説検定の手順

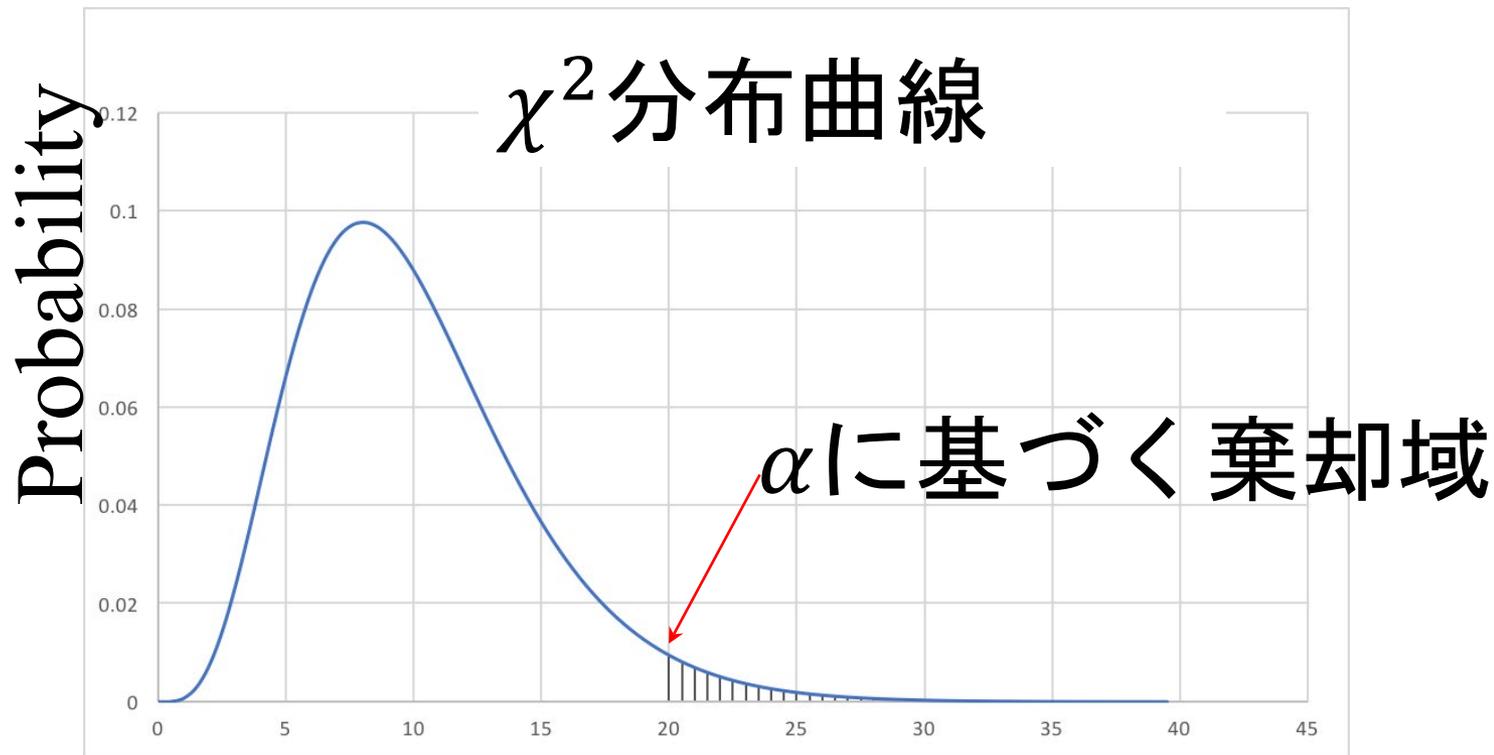
1. 帰無仮説 $H_0$ ， 対立仮説 $H_1$ を決める.
2. 得られたデータから統計量を求める.
  - 用いる統計量：T（T分布）， F（F分布），  
 $\chi^2$ （ $\chi^2$ 分布）
3. 用いる統計量が確率分布にどれだけ従っているかを表す確率 $p$ 値を求める.  $p$ 値は帰無仮説が正しい確率とも言われ， 有意水準 $\alpha$ より小さければ， 帰無仮説 $H_0$ は棄却し対立仮説 $H_1$ を採用する.

# 独立性検定

- 帰無仮説：2変数間が独立
- 対立仮説：2変数間が従属
- 一般的に $\chi^2$  統計量を用いて自由度dfの $\chi^2$ 分布との適合度により独立性を検定する

$$\begin{aligned} \text{検定統計量} &= \sum \sum \frac{(\text{観測度数} - \text{期待度数})^2}{\text{期待度数}} \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{(x_{ij} - y_{ij})^2}{y_{ij}} \end{aligned}$$

# 例：自由度10の $\chi^2$ 検定



- 検定方法：  
p値  $< \alpha$   $\rightarrow$  従属  
p値  $> \alpha$   $\rightarrow$  独立

# 仮説検定法の問題点

- 検定の精度：p値と有意水準  $\alpha$  に依存する.

## これによって引き起こされる問題

- 真に帰無仮説が正しいが，誤って棄却してしまう.  
→ 第一種の過誤 (Type I error) と呼ばれる.
- 真に対立仮説が正しいが，帰無仮説を棄却しない.  
→ 第二種の過誤 (Type II error) と呼ばれる.

# ベイズ的アプローチによる検定

- Bayes factor :  
二つのモデルの周辺尤度の比により検定する.
- 漸近的に真の独立性検定が可能である.
- データセットを $\mathbf{X}$ , 独立なモデルを  
 $g_1: p(x_1, x_2) = p(x_1)p(x_2)$ , 従属なモデルを  
 $g_2: p(x_1, x_2) = p(x_1|x_2)p(x_2)$ としたときの  
周辺尤度の比 Bayes factor (BF) :

$$BF = \frac{p(\mathbf{X} | g_1)}{p(\mathbf{X} | g_2)} \quad \begin{array}{l} BF > 1 : \text{独立} \\ BF < 1 : \text{従属} \end{array} \quad \text{と判定する}$$

# シミュレーション実験1

## -Type I errorの検証-

- 2ノード間が真に独立である構造を用いて実験を行う.
- $\chi^2$  統計量を用いた検定ではデータ数を増やしたとしても Type I errorが発生するが, Bayes factorでは漸近的に収束することを示す.

# 実験結果 - 確率パラメータ: 0.8

## 表: Type I errorの発生率

	10	50	100	500	1000	5000
BF	0.26	0.07	0.02	0.0	0.0	0.0
$\chi^2$	0.16	0.0	0.0	0.03	0.08	0.07
$G^2$	0.17	0.05	0.02	0.03	0.08	0.06

※BF: Bayes factor

# 実験結果 - 確率パラメータ: 0.7

表: Type I errorの発生率

	10	50	100	500	1000	5000
BF	0.23	0.09	0.03	0.02	0.0	0.0
$\chi^2$	0.08	0.08	0.07	0.07	0.05	0.02
$G^2$	0.14	0.11	0.08	0.07	0.05	0.03

※BF: Bayes factor

# 実験結果 - 確率パラメータ: 0.6

表: Type I errorの発生率

	10	50	100	500	1000	5000
BF	0.12	0.04	0.01	0.01	0.0	0.0
$\chi^2$	0.02	0.06	0.04	0.14	0.03	0.07
$G^2$	0.08	0.06	0.04	0.14	0.03	0.06

※BF: Bayes factor

# シミュレーション実験2

## -Type II errorの検証-

- 2ノード間が真に従属である構造を用いて実験を行い, Type II errorの発生率とp値を検証する.

# 実験結果 - 確率パラメータ : 0.8

表 : Type II errorの発生率

	10	20	30	40	50	100
BF	0.3	0.29	0.19	0.11	0.02	0
$\chi^2$	0.61	0.42	0.19	0.1	0.02	0
$G^2$	0.49	0.32	0.18	0.11	0.02	0

※BF: Bayes factor

# 実験結果 - 確率パラメータ: 0.7

表: Type II errorの発生率

	10	20	30	40	50	100	200
BF	0.62	0.61	0.45	0.44	0.31	0.09	0
$\chi^2$	0.89	0.75	0.43	0.39	0.23	0.02	0
$G^2$	0.72	0.65	0.43	0.37	0.24	0.02	0

※BF: Bayes factor

# 実験結果 - 確率パラメータ: 0.6

表: Type II errorの発生率

	40	50	100	200	500	1000
BF	0.77	0.7	0.65	0.33	0.05	0
$\chi^2$	0.73	0.62	0.54	0.19	0.01	0
$G^2$	0.73	0.61	0.54	0.19	0.01	0

※BF: Bayes factor

# 仮説検定の比較

従来の仮説検定では、結果が不安定で必ず Type I 誤差が残るのに対して、ベイズ検定では漸近的に正しい仮説を選ぶことができる。

## 8.まとめ

1. 事前分布は人間の直観では無視されるが数  
学的推論では非常に重要→ないと変な推論  
が起こる
2. ベイズ推定では事後分布の推定がメイン。  
点推定では 事後分布の期待値EAPが用い  
られる。
3. 事後分布の近似的推定にはMCMCもしくは  
変分ベイズが用いられる。
4. モデルを選択するにはモデルの事後確率を  
考える。→周辺尤度
5. 同時確率分布の推定にはAICアプローチ