

# ベイズ統計の導入と行動モデル

2024/09/11 行動モデル夏の学校 @ 本郷

金沢大学 中西航

nakanishi [at] se.kanazawa-u.ac.jp

# ベイズが難しいいくつかの理由

---

- 問題1: ベイズ多義的すぎ
  - 人や文脈によって何を指しているのかが異なる
- 問題2: 講義できちんと扱えない
  - 子供だまし程度しか教える時間がない
- 問題3: 「頻度主義」「ベイズ主義」という空論
  - 50年前に終わったはずの議論がいまだに行われる
- 問題4: 理論がいまだに日進月歩
  - わかっていないことがたくさんある

# 1: いろいろなベイズ

- 歴史上の人物: **Bayes, Thomas (1701?-1761)** イギリスの牧師
- **ベイズ**の定理 by Bayes:  $P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$ 
  - 内容は単なる式変形
- **ベイズ**推定 by Laplace:  $P(X|Y) = \frac{P(X, Y)}{P(Y)}, P(Y|X) = \frac{P(X, Y)}{P(X)}$ 
  - ベイズの定理を用いた確率的推測の方法
- **ベイズ**統計学: より一般の枠組みを指すことば(?)
- 具体的なものには
  - ベイジアンフィルタ
  - ベイジアンネットワーク
  - 状態空間, 隠れマルコフ
  - ...

- サンプルに基づいて(行動)モデルのパラメータを推定する手法のひとつ
  - 「最尤推定」の仲間(対立相手ではない, 単なる代替)
  - もっと言うと「熟練者の直感」「偉い人の決め打ち」なども仲間
- そのような手法の中で, ベイズ統計学に基づいたもの
  - ベイズ統計学に基づく方法=ベイズ推定にしかない優れた点がある
  - もちろん劣る点もある

## 2: 教える時間がない

### ■ 結果的にこういう説明になりがち

- (なんかよくわからんが)パラメータの事前分布  $\varphi(\theta)$
- (最尤推定と同じ)尤度関数  $\prod_i p(X_i | \theta)$
- (掛けて正規化すると)パラメータの事後分布  $p(\theta | X_{1,\dots,N})$ :

$$p(\theta | X_{1,\dots,N}) = \frac{1}{Z} \prod_i p(X_i | \theta) \varphi(\theta) \longleftrightarrow P(Y | X) = \frac{P(X | Y) P(Y)}{P(X)}$$

- 必要ならパラメータの推定量として
  - 事後確率最大化
  - 事後平均

### ■ 当然以下の疑問がある

- 事前分布って何??
- 尤度関数使うなら最尤推定でいいのでは??

## ■ 事前分布って何??

- 分析者の先見的知識を反映(?)
- 階層化して個人差を表現(?)

## ■ 事後分布って何??

- 尤度関数だけでなく事前の情報を反映できるから良い(?)

- 説明すればするほど、そんな主観的でよいの? という批判を際立たせる悪手かもしれない

# でも、便利な例で押し切る

- 階層ベイズあるいはハイパーパラメータ
- サッカー選手の(真の)パス成功率を、ある試合で観測したパス本数・パス試行数から推定したい
  - その日の好不調や対戦相手に依存する
  - フォワードの選手は良い選手でもパス成功率は少ない
- 各選手の真のパス成功率 $q_i$ 、パス試行本数 $N$ のとき、パス成功本数 $k$ となる確率  ${}_N C_k \{q\}^k \{1-q\}^{N-k}$ 
  - $q=0.9$ ,  $N=30$ のとき $k=27$ は24%,  $k=24$ は5%
  - たとえば $k=24$ なら $q$ の最尤推定量は0.8だが、何人もいたら真の $q$ と遠い観測量はたくさん得られるはず

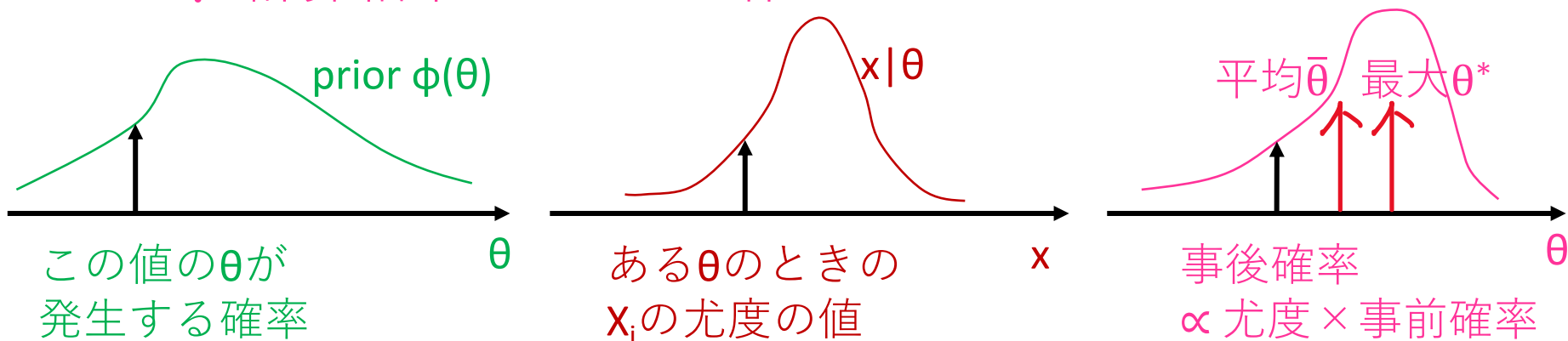
# でも、便利な例で押し切る

- 各選手の $k_i/N_i$ だけから $q_i$ をもう少しよく推定したい
  - 事前分布 $\varphi(q)$ : 試合に出てる選手はある程度上手い**と思う**ので、たとえば $\varphi(q) \sim N(0.9, 0.01)$ などとする
  - もう少しきちんと $q_i = a_i + b_i$ などを考えることもできる
  - $a$ や $b$ が従う分布を考える: たとえば,
    - $a_i$ は全選手の平均(共通の定数)とする
    - $b_i$ は正規分布 $N(0, \sigma^2)$ に従うとする
  - $\sigma$ の事前分布を考える
    - 階層化/ハイパーパラメータと呼ばれるもの

よく見かけるベイズ推定の1つ目の顔:  
**異質性の表現としてのベイズ推定**

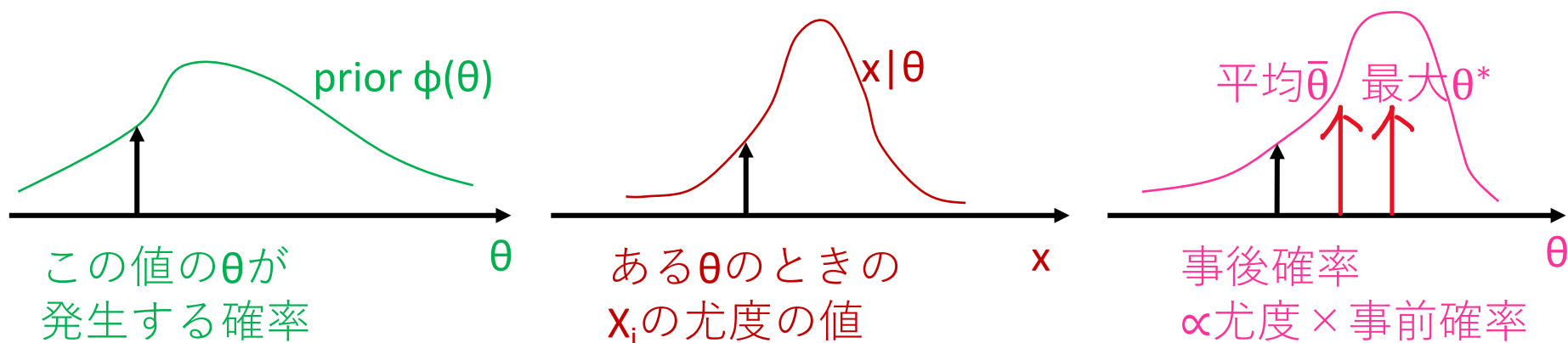


- どうやったら計算できる??  $p(\theta | X_{1,\dots,N}) = \frac{1}{Z} \prod_i p(X_i | \theta) \varphi(\theta)$
- 解析的に行えるケースは稀(共役事前分布)
- 一般には数値的に近似する
  1. 事前分布に従って $\theta$ をひとつ選ぶ
  2. その $\theta$ で尤度を計算する
  3. 計算結果をたくさん作る



- $Z$ を求めるのは難しいが、 $Z$ がなくても事後分布の形はわかる
- 最大や平均もわかる

- 数値的に近似する方法の別の見方...



- もし、事後分布の平均値を与える $\bar{\theta}$ を得たいだけなら、良い計算の方法がある(説明省略)
- 仮に尤度関数の形が複雑でも、**尤度関数を直接最大化しなくても、 $z$ を求めなくても**、たくさん計算すれば関数の形がだいたい求まる

よく見かけるベイズ推定の2つ目の顔:

**アルゴリズムとしてのベイズ推定**

# 3: 確率的主義という虚構

---

説明の前に...

- これから、よく聞くであろう、もしかすると講義で習うであろうことはおよそ間違いだ、という話をします
- このことは統計学・数理科学の人々は当たり前のように認識していることです
- が、それ以外の分野ではいまだに根強い誤解です

## ■ モデル推定(統計的推測・学習)

- 得られているサンプル(データ)に基づいて、  
真の確率分布(モデル)に接近しようとする試み

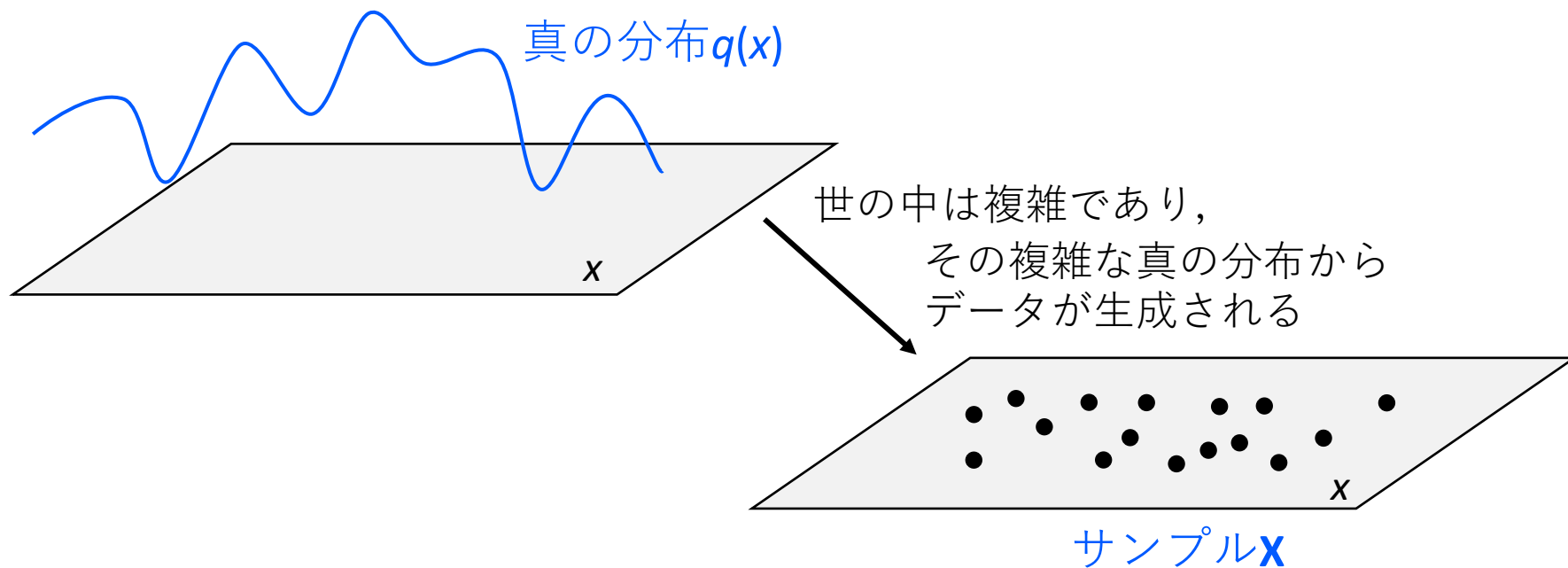
## ■ モデル評価

- 推定したモデルの良さを客観的・相対的に比較すること

## ■ モデル選択

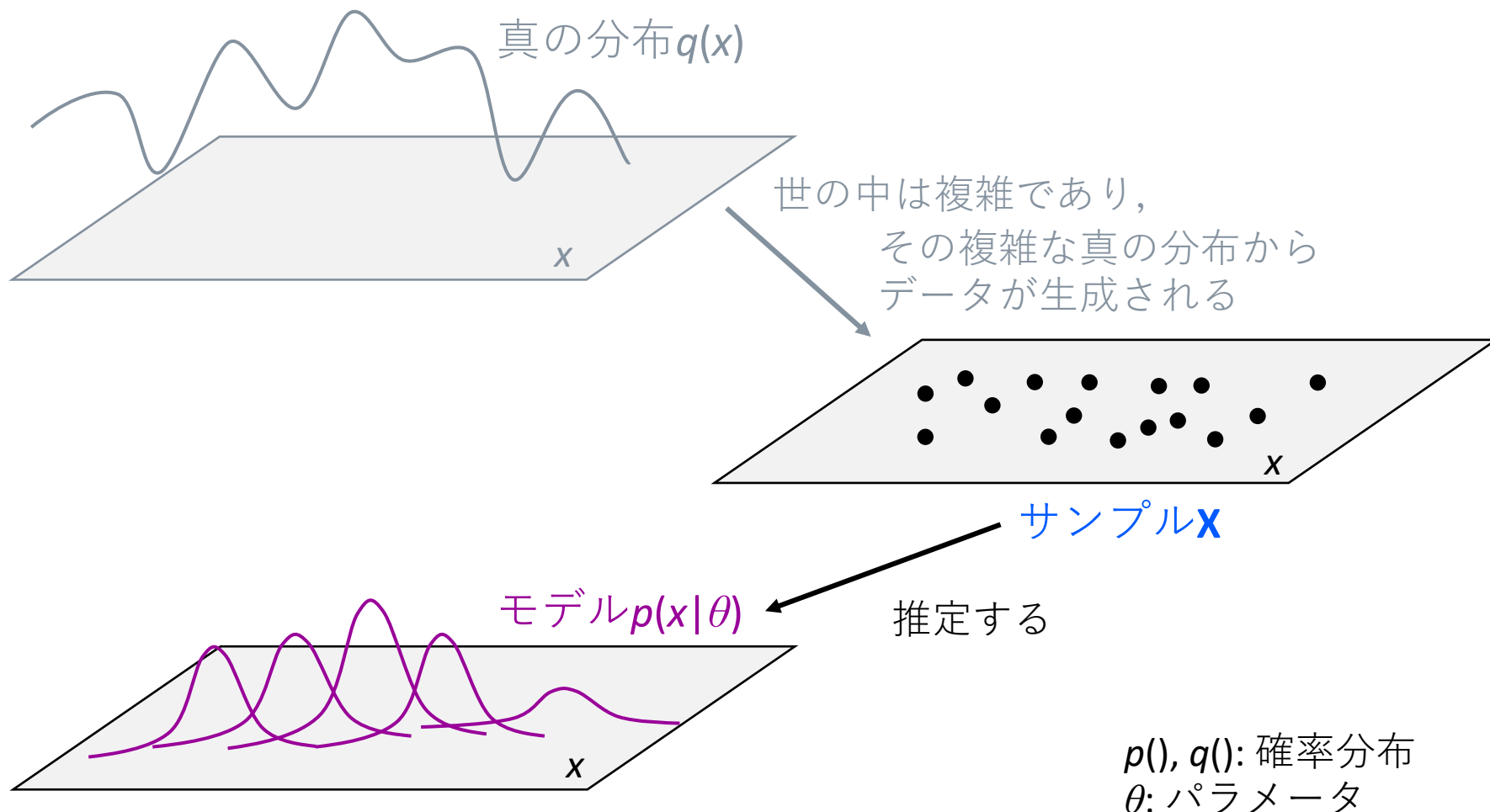
- 考えうるモデル候補から良いモデルを選択すること
- “良い”の意味は後述

- 真の分布 $q(x)$ (=モデル)をサンプル $\mathbf{X}$ (=データ)から推定したい

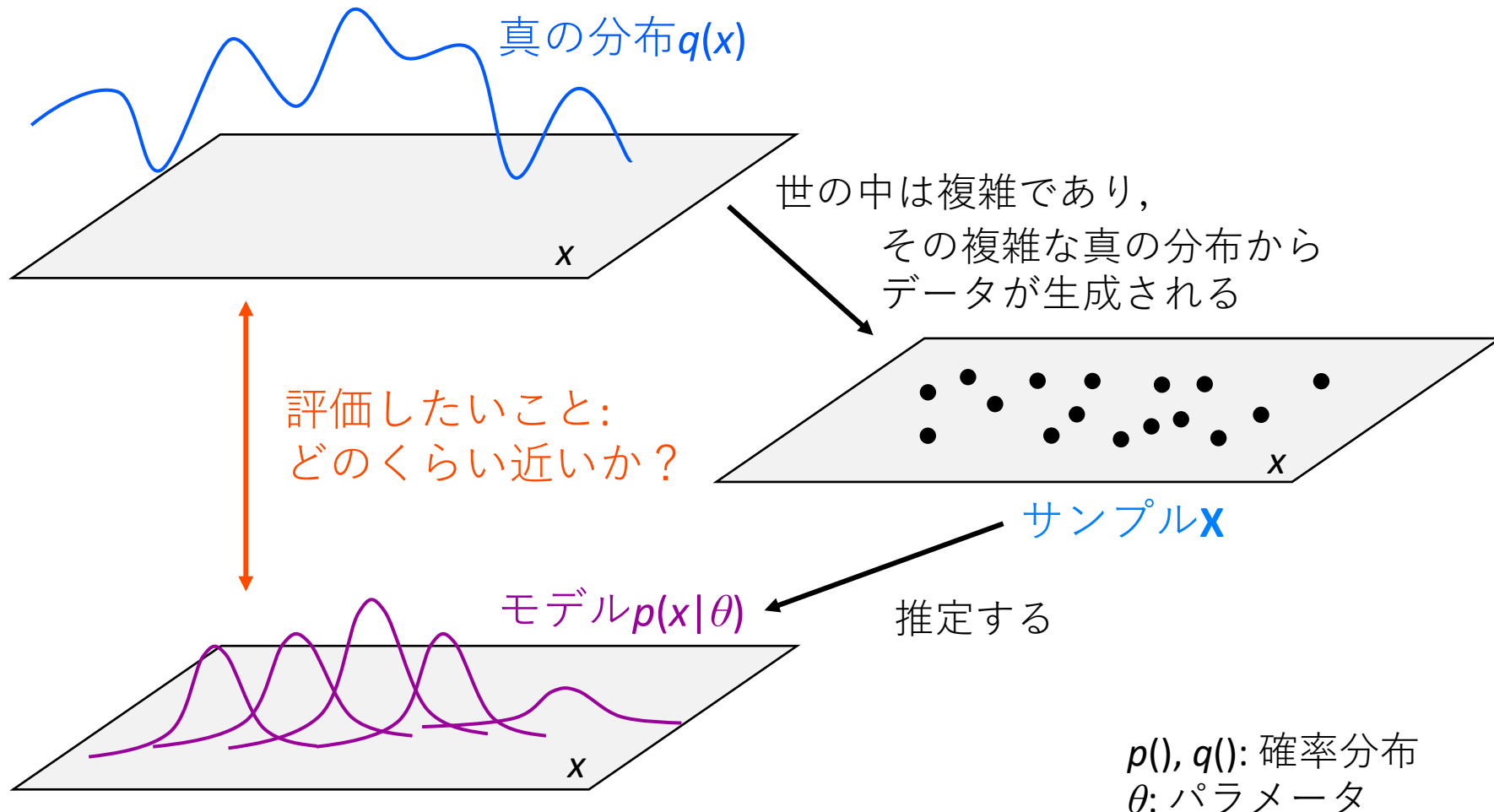


$q(\cdot)$ : 確率分布

- 真の分布 $q(x)$ (=モデル)をサンプル $\mathbf{X}$ (=データ)から推定したい



- 真の分布 $q(x)$ (=モデル)をサンプル $\mathbf{X}$ (=データ)から推定したい



- モデル推定:  $q(x)$  になるべく近い  $p(x|\theta)$  を求めること
- モデル評価:  $q(x)$  に  $p(x|\theta)$  がどの程度近いかを見積もること
  - 近いほど良いモデルと考える
- 統計的モデル評価の本質的な意味:
  - 真の分布  $q(x)$  が未知でも, サンプル  $\mathbf{X}$  のみから,
  - (1) モデル  $p(x|\theta)$  と真の分布  $q(x)$  との近さがわかる
  - (2) モデル  $p(x|\theta)$  による未知データの予測精度がわかる
    - より正確には, それらの期待値が分かる
    - そうでなければ意味がない(手元のデータで評価不可能)
- これはビッグデータでもスモールデータでも一緒



## ■ モデル候補(モデルA, モデルB, ...)のなかで

- (1) 真の分布 $q(x)$ に最も近いモデル
- (2) 未知データの予測精度が最も良いモデル

## ■ 真の分布がモデルに対して**正則な場合**(条件は省略),

- (1) は Bayesian Information Criterion (BIC) 最小のモデル
- (2) は Akaike Information Criterion (AIC) 最小のモデル

## ■ この2つは一般に一致しない

※ 厳密には、AICは実現可能なときに限る(説明省略)

- 善し悪しではなく、分析者が何を求めるか

## ■ これらが唯一(唯二?)の方法ではない

- 尤度も、統計的検定(尤度比)も、信念もダメではない
- どのような根拠に基づくかの認識と使い分けが重要

# 確率論と主義: モデル評価には無関係<sup>19</sup>

- このような整理を見ることがあるが...

	頻度主義的	ベイズ的
真のモデル	神のみ 絶対的に正しい単一のもの	分布 $q(x)$
サンプル $x$	頻度主義的な真のモデル + なんらかの分布に従う誤差	ベイズ的な真のモデルの サンプル
推定モデル	揺らがない $f(x)$ とサンプルに よって発生する誤差分布	サンプルによって揺らぐ $q(x)$

こんなものはない!!!

- $f(x)$  + (誤差分布) を  $q(x)$  とみれば両者はまったく一緒
  - $f(x)$  が揺らいでいないというのは信念であり識別不可能
- 主義・信念とモデルの推定方法(最尤/ベイズ)は無関係
- 主義・信念はモデルの評価・選択においては無意味

-1930

1950

1970

1990

2010-

## 主義によるモデル選択

## 規準によるモデル選択

最尤法 (Fisher 1912-22)

AIC (Akaike 1974)

WAIC (Watanabe 2009)

ベイズ法 (もっと昔から)

BIC (Schwarz 1978)

WBIC (Watanabe 2013)

### 正則な場合

### 一般の場合

- 真の分布が未知だから  
正確さは永遠にわからない
- 確率とはなにかという決め事  
確率についての哲学の違い
- 最尤法は客観的  
ベイズ法は主観的
- 主義主張の対立と停滞

- 真の分布が未知でも  
正確さを見積もれる
- 統計的モデル推定には無関係  
哲学ではなく数理科学の問題
- どちらにせよモデルは主観  
ゆえに客観的な評価規準が重要
- パラダイムシフトによる進展

■ モデルの統計的な側面は、工学ではなく理学の興味の範疇と捉えられることがあるが、

■ むしろ、やっと工学で議論できる基盤ができてきた

	最尤推定	ベイズ推定
目的	サンプル $\mathbf{x}$ から真の分布 $q$ に接近したい	
用意するもの	サンプル $\mathbf{x}$ とモデル $p(\mathbf{x} \boldsymbol{\theta})$ ( $\rightarrow$ 尤度関数)	
仮定	<b>尤度最大が最良</b>	<b>事前分布<math>\varphi(\boldsymbol{\theta})</math>の存在</b>
得られる 推定量・分布	最尤推定量 $\hat{\boldsymbol{\theta}}$ 必要なら, 分散共分散行列 (=尤度関数の形状)	事後分布 $p(\boldsymbol{\theta} \mathbf{x})$ 必要なら, 事後確率最大化・ 事後平均等の代表値
得られる モデル	最尤推定量を モデルに代入したもの	事後分布で モデルを平均したもの
使えるとき	<b>正則な場合のみ</b>	<b>一般の場合</b>

- 「最尤推定は点推定、ベイズ推定は区間推定」  
→ 誤りではないが、それはこの部分だけ

- 事前分布の存在を認めるor 尤度最大が最良と信じるという選択をしている

目的		
用意するもの		
仮定	<b>尤度最大が最良</b>	<b>事前分布<math>\varphi(\theta)</math>の存在</b>
得られる推定量・分布	最尤推定量 $\hat{\theta}$ 必要なら、分散共分散行列 (=尤度関数の形状)	事後分布 $p(\theta X)$ 必要なら、事後確率最大化・事後平均等の代表値
得られるモデル	最尤推定量をモデルに代入したもの	事後分布でモデルを平均したもの
使えるとき	<b>正則な場合のみ</b>	<b>一般の場合</b>

- 最尤推定量の最良性は正則な場合のみ成立

- 最尤推定も真のモデルに近いモデルを推定する作業
- 最尤推定量はその過程で出てくるもの

- パラメータの事後分布と最終的に求まるモデルは別物

モデル  $V = \beta_0 + \beta_1 x + \epsilon, \epsilon \sim (\text{何らかの確率分布})$

パラメータ  $\beta_0, \beta_1$  を推定し, モデルを推定する

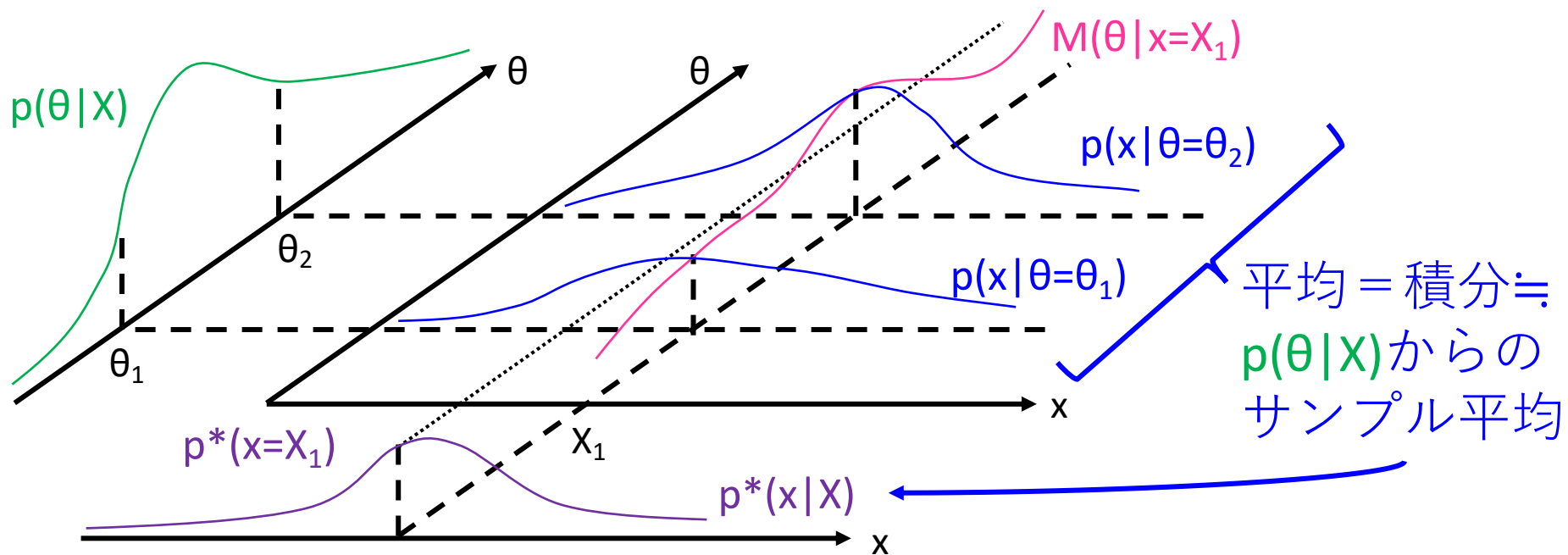
## ■ 最尤推定

- 最尤推定量  $\widehat{\beta}_0, \widehat{\beta}_1$
- それを用いたモデルは  $V^* = \widehat{\beta}_0 + \widehat{\beta}_1 x + \epsilon$

## ■ ベイズ推定

- 事前分布  $\varphi(\beta_0, \beta_1)$
- 事後分布  $p(\beta_0, \beta_1 | \mathbf{X})$
- それを用いたモデルは  $V^{**} = \iint p(V)p(\beta_0, \beta_1 | \mathbf{X})d\beta_0d\beta_1$

## ■ モデルを事後分布で平均



推定された(真の分布に最接近した)モデル

# (行動)モデルのベイズ推定



- 顔2: アルゴリズムとしてのベイズ推定が発端
  - Probit等のopen formの推定を行う際の手段  
= 尤度関数の最大化が難しい場合
- たとえばTrain本の12章はほとんどこの話

	Supplier	270
11.7	Discussion	280
12	Bayesian Procedures	282
12.1	Introduction	282
12.2	Overview of Bayesian Concepts	284
12.3	Simulation of the Posterior Mean	291
12.4	Drawing from the Posterior	293
12.5	Posteriors for the Mean and Variance of a Normal Distribution	294
12.6	Hierarchical Bayes for Mixed Logit	299
12.7	Case Study: Choice of Energy Supplier	305
12.8	Bayesian Procedures for Probit Models	313
13	Endogeneity	315
13.1	Overview	315
13.2	The BLP Approach	318

- 顔2: アルゴリズムとしてのベイズ推定が発端
  - Probit等のopen formの推定を行う際の手段  
= 尤度関数の最大化が難しい場合
- Train本, 12.1

## 12.1 Introduction

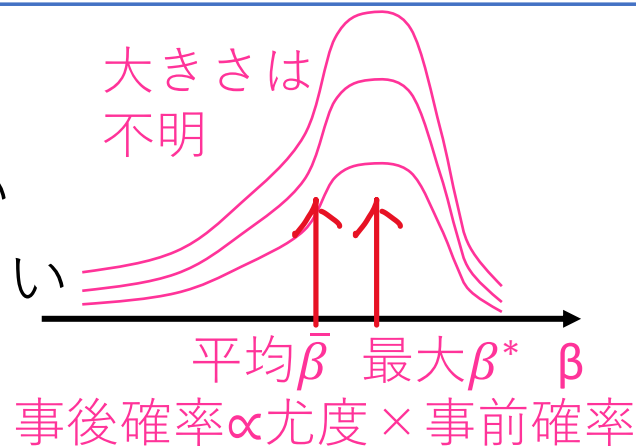
A powerful set of procedures for estimating discrete choice models has been developed within the Bayesian tradition. The breakthrough concepts were introduced by Albert and Chib (1993) and McCulloch and Rossi (1994) in the context of probit, and by Allenby and Lenk (1994) and Allenby (1997) for mixed logits with normally distributed coefficients. These authors showed how the parameters of the model can be estimated without needing to calculate the choice probabilities. Their procedures provide an alternative to the classical estimation methods described in Ch

- 尤度の計算には選択確率の計算が必要だが、open formでは難しい
- ベイズ推定でこれを回避

- 顔1: 異質性の表現としてのベイズ推定の話もある
- Train本, 12.1, つづき

described in Chapter 10. Rossi *et al.* (1996), Allenby (1997), and Allenby and Rossi (1999) showed how the procedures can also be used to obtain information on individual-level parameters within a model with random taste variation. By this means, they provide a Bayesian analog to the classical procedures that we describe in Chapter 11. Variations of these procedures to accommodate other aspects of behavior have been numerous. For example, Arora *et al.* (1998) generalized the mixed logit procedure to take account of the quantity of purchases as well as brand choice in each purchase occasion. Bradlow and Fader (2001) showed how similar methods can be used to examine rankings data at an aggregate level rather than choice data at the individual level. Chib and Greenberg (1998) and Wang *et al.* (2002) developed methods for interrelated discrete responses. Chiang *et al.* (1999) examined situations where the choice set that the decision maker considers is unknown to the researcher. Train (2001) extended the Bayesian procedure for mixed logit to nonnormal distributions of coefficients, including lognormal, uniform, and triangular distributions.

- $Z$ を求めないと形しかわからないので、モデル評価のための統計量は求まらない  
→ベイズ推定の**統計的**利点は発揮されない



- 本来ベイズ推定で得られるモデルは  $V^{**} = \iint p(V)p(\beta_0, \beta_1|\mathbf{X})d\beta_0d\beta_1$  だが、この使い方は見たことがない。
- むしろ、事後確率の平均を与える $\bar{\beta}$ ・事後確率の最大を与える $\beta^*$ を代入したモデルをよく見かける： $V^{***} = \bar{\beta}_0 + \bar{\beta}_1x + \epsilon$ 
  - これは、厳密にはベイズ推定の結果のモデルではない
  - 正確には「最尤推定の推定プロセスをベイズ推定に、尤度最大を事後確率の平均で置き換えたもの」
  - 世の中ではこれらを「ベイズ推定」と称することも多い

Two important notes are required before proceeding. First, the **Bayesian procedures**, and the term “hierarchical Bayes” that is often used in the context of discrete choice models, **refer to an estimation method, not a behavioral model**. Probit, mixed logit, or **any other model** that the researcher specifies can, in principle, **be estimated by either classical or Bayesian procedures**. Second, the Bayesian perspective from which these procedures arise provides a rich and intellectually satisfying

- 「ここでいう Bayes とは、行動モデルが Bayes というのではなく、推定方法が Bayes ということ」

- 「分析者はあらゆるモデルを Classical な方法でも Bayes の方法でも推定できる」

■ ただし、統計的にきちんと説明するならば、

- “(行動)モデルが Bayes か否か”は概念であり、統計的な区別ではない
- ここでいう“Bayes の方法”はあくまでも尤度関数の何らかの意味での最大化



# 行動モデル in ベイズ的なモデルの例 32

- 観測された一連のGNSS座標( $y$ )から、直接観測のできない歩行者の存在リンクおよび経路( $x$ )を推定する問題

Literature review Introduction

## Route measurement models (2)

*Path-based probabilistic approach* evaluates path likelihood regarding all GPS data included in a trip  
Pyo et al. (2001); Bierlaire et al. (2013)

Domain of Data Relevance

Likelihood

Path 1: 10%

Path 2: 40%

Path 3: 50%

Oyama, Y. (The University of Tokyo) 15<sup>th</sup> Behavior Modeling summer course Sep. 25, 2016 9

Literature review Introduction

## Route measurement models (3)

*Bayesian approach* incorporates behavioral models into measurement models  
Chen et al. (2013); Danalet et al. (2014)

Prior (Route choice model)

Measurement

Posterior

Oyama, Y. (The University of Tokyo) 15<sup>th</sup> Behavior Modeling summer course Sep. 25, 2016 12

<http://bin.t.u-tokyo.ac.jp/model16/lecture/Oyama.pdf>

- (注) 大山さんの上記スライドはあくまで研究背景で、Oyama and Hato (2018)は上記手法に内包される問題を解決した研究

## ■ [正当な批判]

- 計算時間がかかりすぎる，いくら計算機が進化したといっても無理なものは無理
- コーディングが面倒で複雑，現場担当者が使えない
- 収束判断がよく分からない
- (?)私のモデルは正則だと信じている．最尤推定で十分

## ■ [誤った批判]

- 事前確率は主観が入っており客観的ではない! 最尤推定こそが客観的である!!!
- 推定量が不偏性を持たないなんてけしからん!!!



## ■ 必要なときはベイズ推定を使いましょう

- 階層化したい，尤度関数の最大化が難しい→もちろんOK
  - ただし，ベイズ統計の本質とはあまり関係ない
- 事前分布への批判は気にしなくてOK
  - モデルの想定が支配的に重要
- 本当は，想定したモデルにおいてベイズ推定が望ましい (=正則でない)ならば，ベイズ推定すべき

## ■ 一方で，どのようなときにモデルが正則でないかは数学の問題で，簡単ではない

- 端的には「ほとんどの場合」
- とはいえ，工学的にはそれだと困る
- 妥協のしどころはよく分かっていない(分野固有の問題)