

# 簡便エミュレーションによる 実験計画の高速化

樋口知之

情報・システム研究機構 統計数理研究所

1/42

## アウトライン

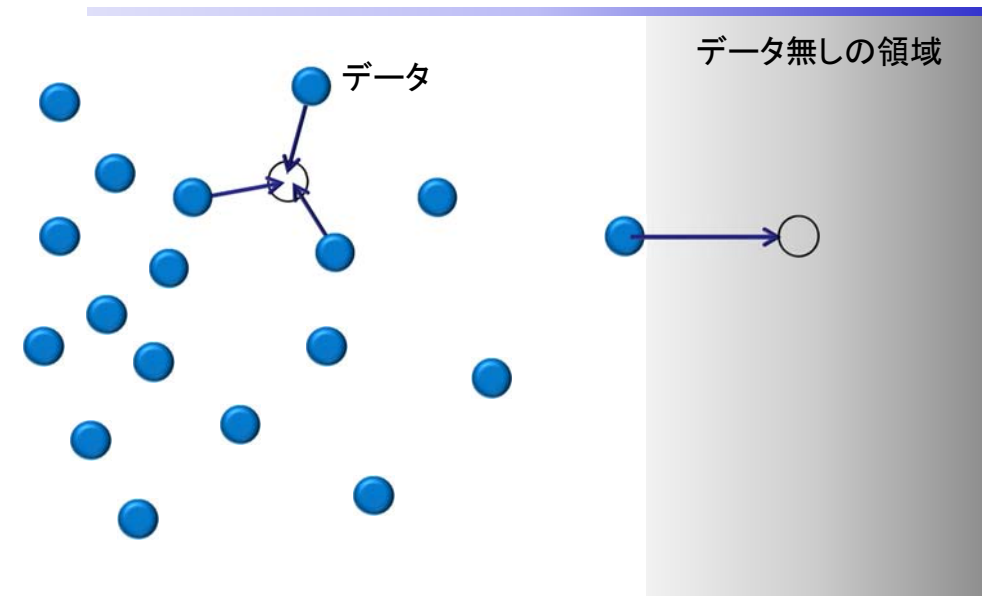
1. データ同化
  - a. 基本概念と目的
  - b. アルゴリズムの基本
  - c. 設計技術の革新
  - d. DAセンターによる応用事例
2. 簡便エミュレータ
  - a. 必要な理由と動向
  - b. 構成法
3. スパース回帰: LASSO, CS, NMF
4. GPR: Gaussian Process Regression

2/42

### 1.a 基本概念と目的

3/42

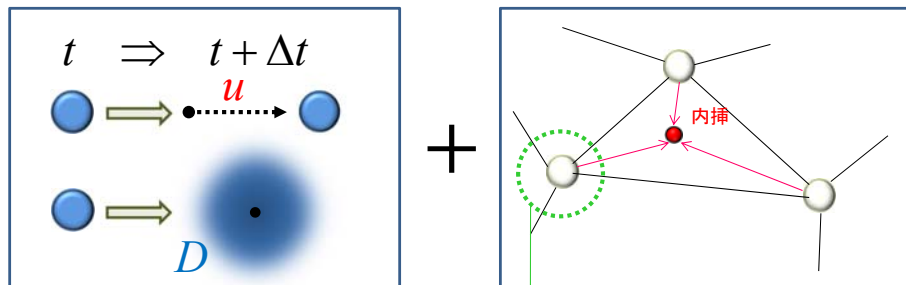
## 内挿と外挿問題



4/42

# いろいろなビジネス展開が可能

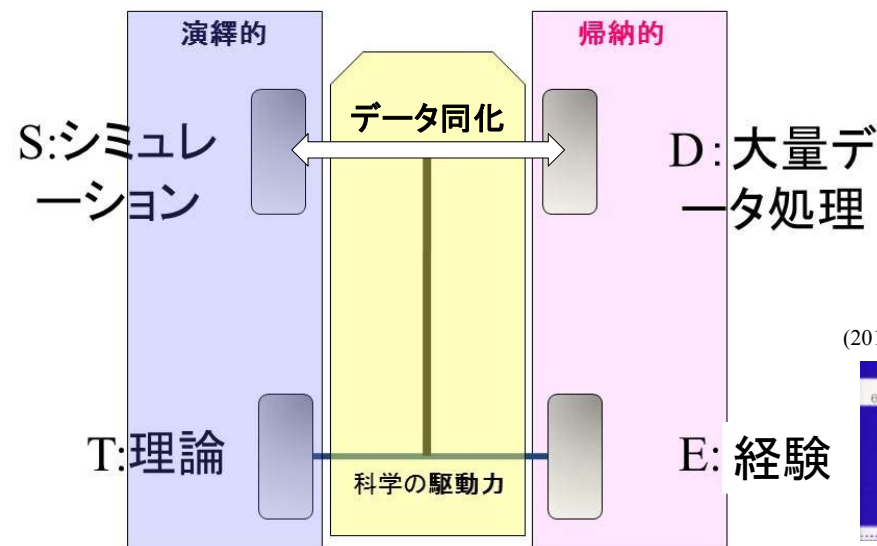
移流と拡散 + クラウドソーシング = 予測能力  
 フォワード計算モデル 現況を捉える認識力 スマホ(とGPS)



情報の不確実性  
 (多人数からのレポート: 多様なノイズの混在)

5/42

# つなぐ：データ同化



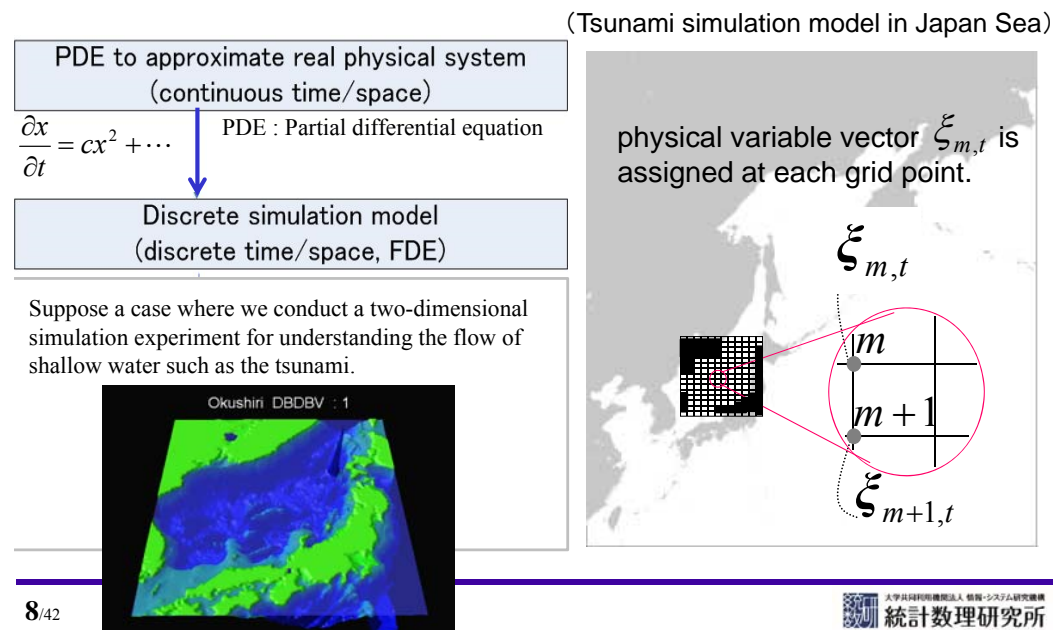
(2011年9月刊行)



6/42

## 1.b アルゴリズムの基本

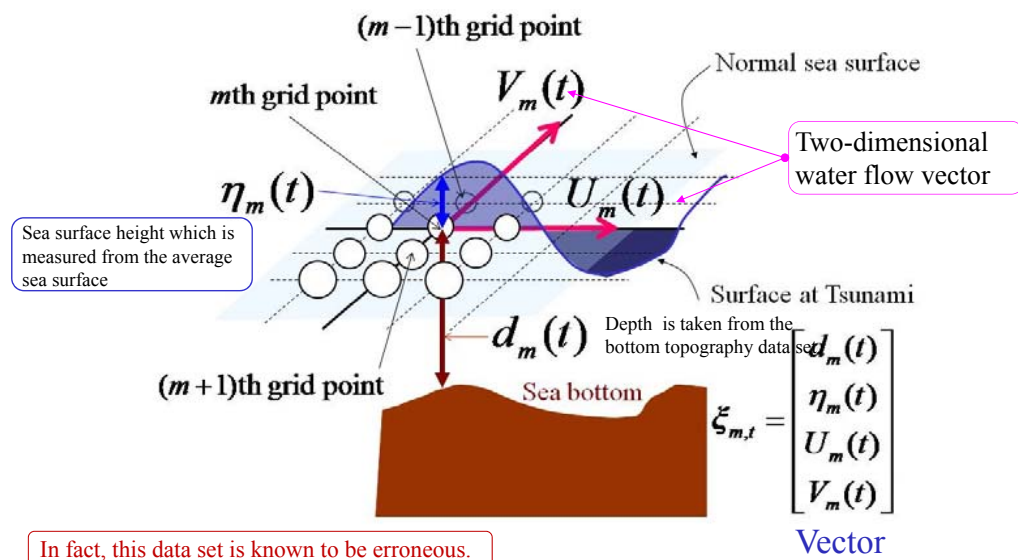
## シミュレーションモデルの構成 (1)



7/42

8/42

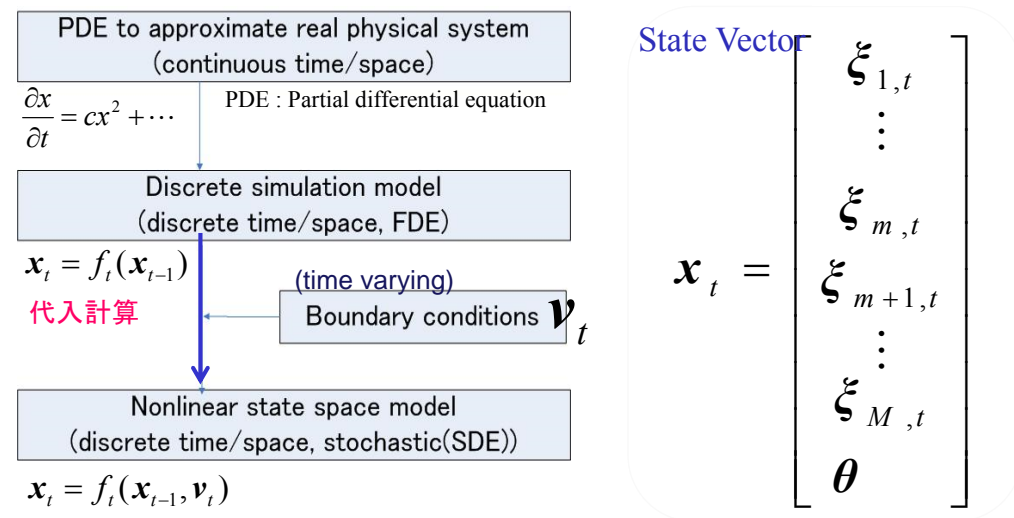
## シミュレーションモデルの構成 (2)



9/42

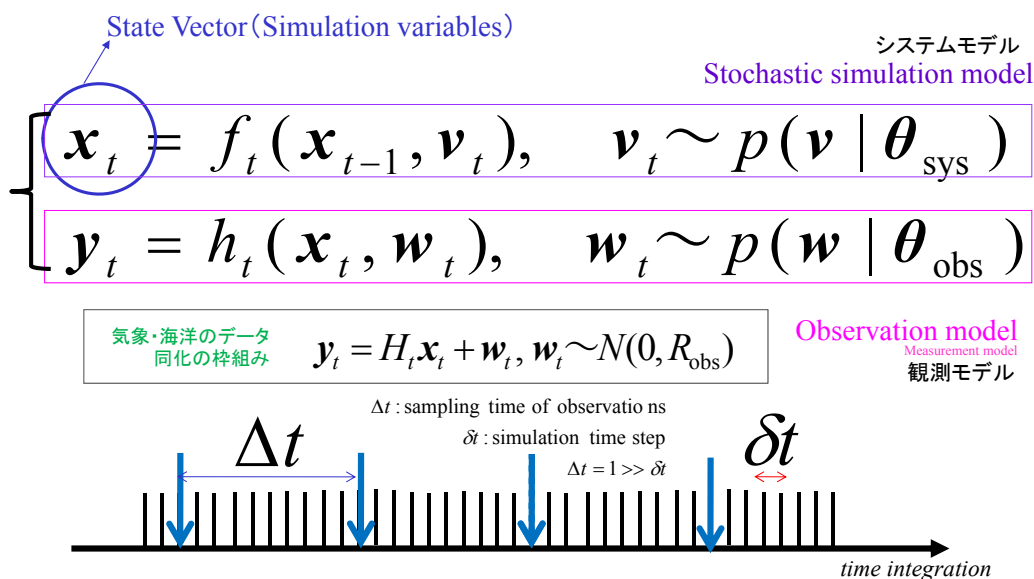
## システムモデルとしてのシミュレーションモデル

(simplified meteorological model around Japan)



10/42

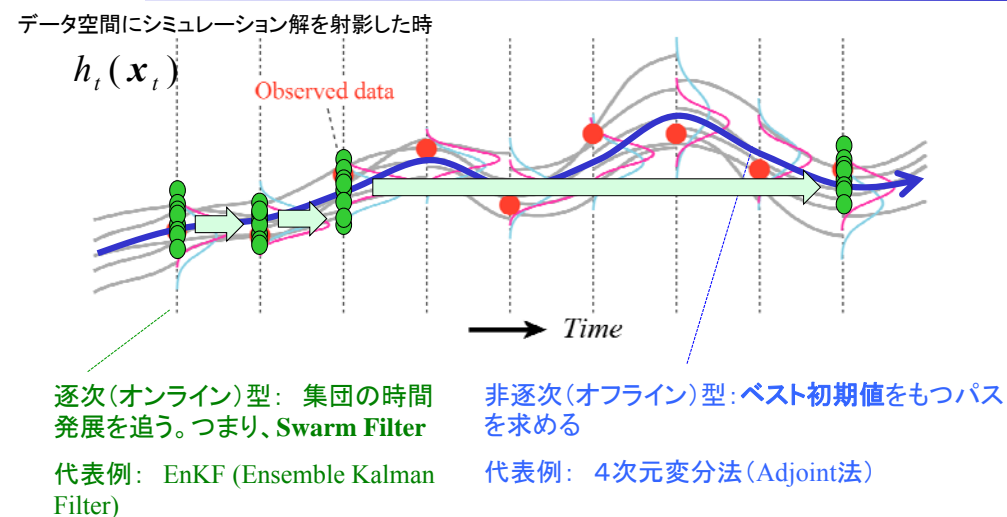
## データ同化と一般状態空間モデル



11/42

## 逐次 (アンサンブル) vs. 非逐次 (最適パス)

データ同化のイメージ



12/42



# マルチエージェントシミュレーションへの応用

Cyber-physical systems (CPS) are engineered systems that are built from and depend upon the synergy of **computational and physical components**.

The term "cyber-physical systems" refers to the tight conjoining of and coordination between computational and physical resources.



An embedding of the multi-agent based mode, that is a person-based model to describe each behavior, can be easily achieved by replacing a **grid** with a **person** in the definition of the state vector.

## Human Modeling

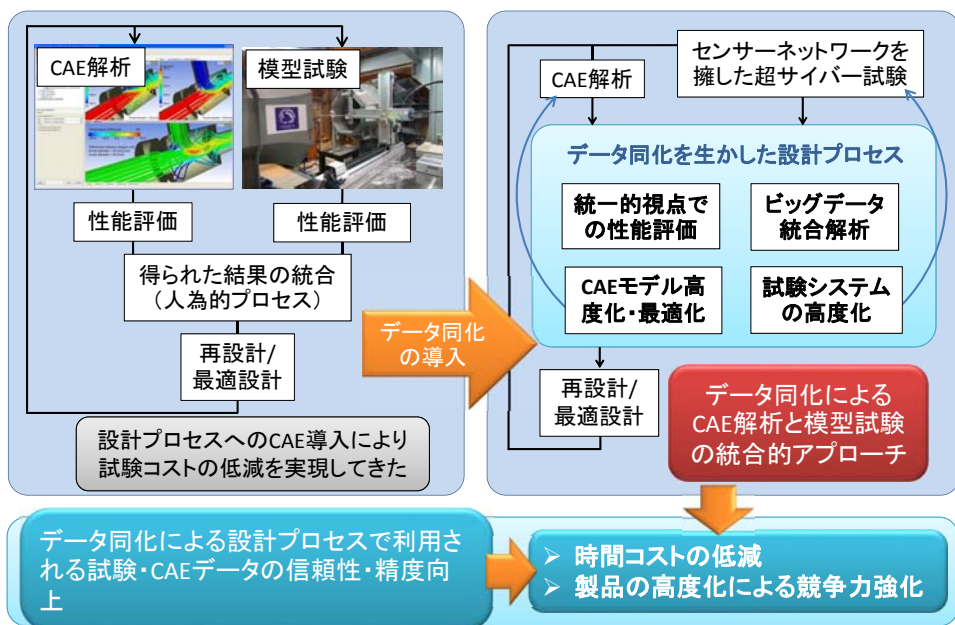
$\xi_{m,t}$  = (variables to specify a behavior of the  $m$ -th agent)

$p(x_t | x_{t-1})$  **MAS: Multi-Agent based Simulation**

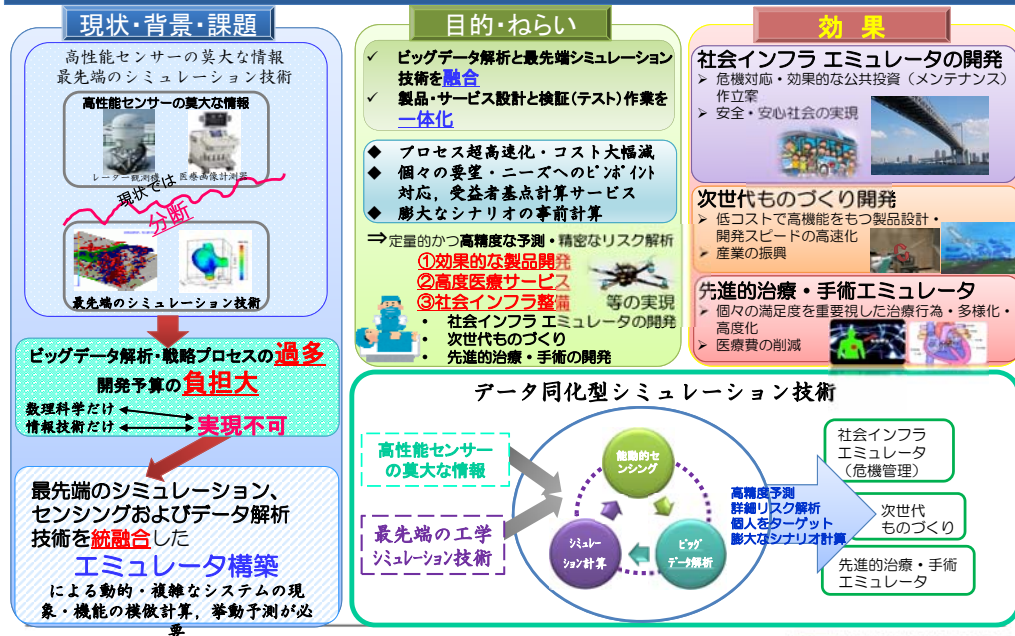
$$x_t = \begin{bmatrix} \xi_{1,t} \\ \vdots \\ \xi_{m,t} \\ \xi_{m+1,t} \\ \vdots \\ \xi_{M,t} \\ \theta \end{bmatrix} \leftarrow f_t \begin{bmatrix} \xi_{1,t-1} \\ \vdots \\ \xi_{m,t-1} \\ \xi_{m+1,t-1} \\ \vdots \\ \xi_{M,t-1} \\ \theta \end{bmatrix}$$

## 1.c 設計技術の革新

## データ同化による設計技術の革新



## データ同化型シミュレーション技術開発によるものづくり・設計の革新 ーデータ同化技術の工学への応用研究開発ー(1/2)





**COCN**  
Council on Competitiveness-Nippon  
HPC応用研究会提言  
産業競争力懇談会  
(2012年3月6日)  
より抜粋・要約

【効果あるモデリング方法論】  
(ア)自然現象を数学モデルに近似するモデルによる誤差 (イ)材料データなど入力データを持つばらつき(構成式、減衰率、熱伝達率など)に対応するシミュレーション技術としては、**DA(データ同化)を推進したい**。DAは複雑な現象を科学的に理解し、精度よく予測したいという要求にこたえる技術であり、モデルから複雑現象を再現する演算アプローチと複雑現象の観測結果からモデルを推測する帰納アプローチが融合した**次世代のシミュレーション技術**である。これまで主に**地球科学の分野**において数値モデルの再現性を高めるためにモデルに観測データを埋め込み、馴染ませることを意図して研究された。**モノづくり分野**において、なじみが薄い名前ではあるが、たとえば、電子機器の設計分野では20年も前から実質的に同等の技術が活用されている。

ー現代社会に強く求められるデータ同化技術の研究開発ー

**社会インフラ  
エミュレータ研究事例**

地盤工学におけるデータ同化  
構造シミュレーション

地盤挙動のシミュレーション解析と将来予測システムの開発  
地盤の変位や応力、間隙水圧の確率分布を得てこれらの確率分布を基礎や土構造物のリスク評価に適用

**地盤内応力状態把握  
精度検証の実現**

**次世代ものづくり開発  
研究事例**

後方乱気流のライダー計測融合  
シミュレーション

次世代の飛行機運航システムの開発  
仙台空港における実運航機の後方乱気流をライダー計測と3次元の流体シミュレーションに融合

**離発着間隔の制限解除  
安全時短省エネ運航の実現**

**先進的治療・  
手術エミュレータ研究事例**

データ同化型血流シミュレーション

超音波計測融合血流解析システムの開発  
血流の超音波計測と数値解析を一体化し血流解析システムを開発し、疾患との関係や診断指標について研究

**高度診断の実現**

## 1.d DAセンターによる応用事例 (ちょっと脇道)

## Application studies carried out by our group for model improvements

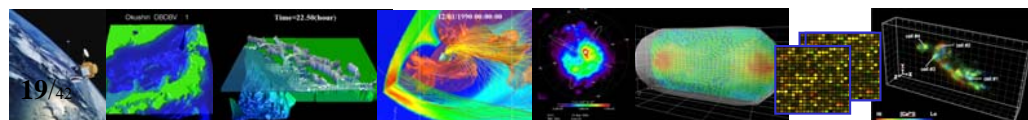
### Research Projects:



- Typhoon trajectory
- Tsunami, Ocean tide
- Auroral phenomena
- 3D structure of ring current
- Acoustic waves
- Intercellular fluid dynamics
- Genome informatics
- Drug response prediction
- Neuronal circuits of whole nerve cells
- Influenza Pandemic (Multi agent simulation)

Watch You Tube "Data Assimilation  
R&D Center at ISM".

さらに モデル性能のヒントおよび研究成果に



## 簡便エミュレータ 2.a 必要な理由と動向

# 通常のエミュレータの概念

エミュレータ (Emulator)とは、コンピュータや機械の模倣装置あるいは模倣ソフトウェアのことである。

## 概要

コンピュータ分野で使われることが多い用語だが、もともとは機械装置全般に使う言葉である。判りやすく言えば、機械を真似る機械である。

## 語源

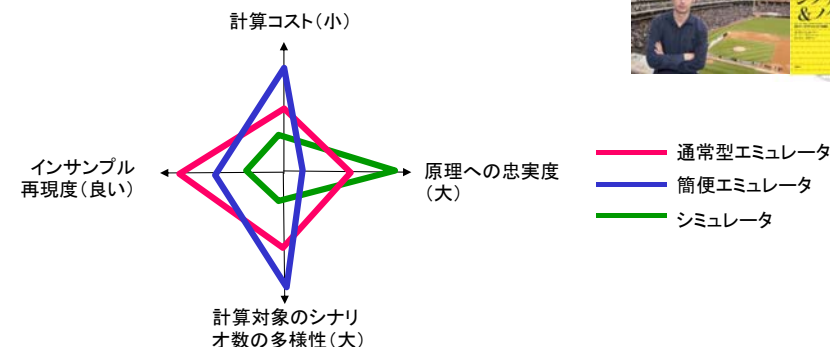
エミュレーションやエミュレータは、模倣対象のシステムにおいて、**予測できる現象より予測できない現象が支配的**である場合に使われる。また、非常に高い安全性が要求される場合にも良く使われる。**予測できる現象が支配的な場合や、完全に模倣することが難しい場合はシミュレーション技術**を使う。

(Wikipediaより)

# 計算コストと予測精度

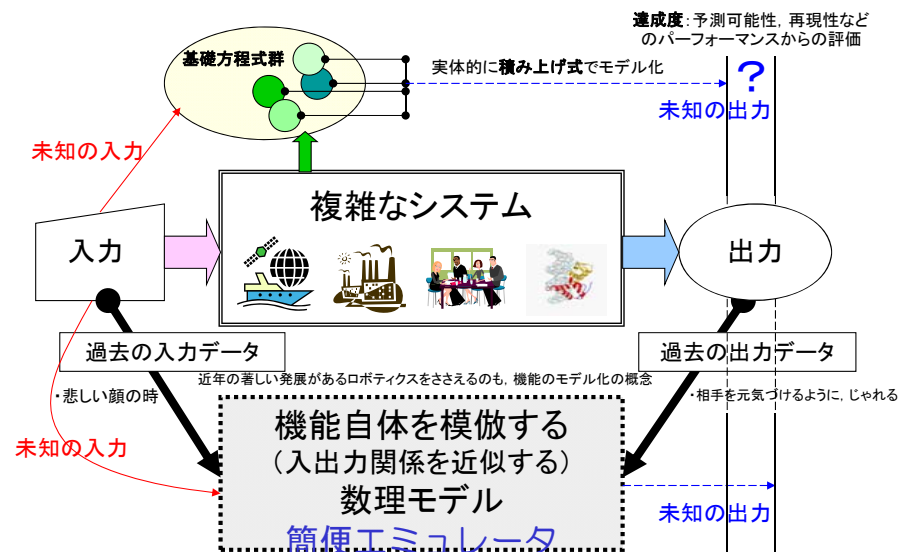
- 高い計算コスト(オンライン計算が難しい)
- 不確実性の取り扱い方

“不確実性を、世界を理解する人間の能力に付随するものではなく、実験に付随するものとしてとらえている。”

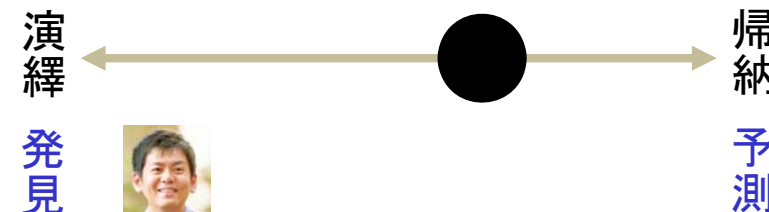


# 機能のモデル化：簡便エミュレータ

高度情報社会におけるユニバーサルな研究課題の表現形



# 予測と発見



鹿島・京大教授

グレイボックスモデル  
(元トヨタ 大嶋氏)



統計数理(2006)  
第54巻 第2号 209-210  
©2006 統計数理研究所

特集「予測と発見」

2006年

「特集 予測と発見」について

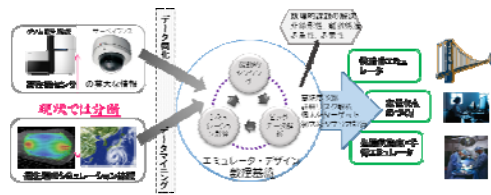
# UQ: Uncertainty Quantification

- 欧米では、計算機シミュレーション結果の信頼性を具体的に確立するための方法論の研究が急速に熱を帯びてきており、ASME(The American Society of Mechanical Engineers)がVerification and Validation (通常V&Vと呼称)の標準化に大きな力を注いでいる。例えば、2006年には固体力学に対して、2009年には流体力学および熱解析に関する計算機シミュレーションのV&Vが公表されている。
- 欧州においては流体力学の分野で同種の研究活動が2012年から活発化しており、Uncertainty Quantification (UQ) in Industrial Analysis and Design の名のプロジェクト研究が現在進行中である。
- NASAでは、NASA UQ challenge 2014と題して、スパースな限定されたパラメータセットに関するシミュレーションの結果データから、UQをモデル化するコンペを開始した。
- 米国統計コミュニティは、2011-12年に、NSFのサポートを受ける機関SAMS(I(Statistical and Applied Mathematical Sciences Institute)にてUQを集中的に研究するプログラムを立ち上げた。
- 米国統計学会はSIAM(Society for Industrial and Applied Mathematics)と共同でJournal on UQの刊行を2014年に開始した。その雑誌の取り扱う主たる分野としてsensitivity analysis, model validation, model calibration, data assimilationの4つがあげられている。最新号の論文(4本掲載)は、感度解析、ガウス過程回帰、モデル較正、ギブスサンプラーの解析のテーマとなっており、ほぼ統計学の範疇である。

重要な技術:

ガウス過程回帰や、その古典版とも言えるクリギング  
次元削減を目的としたスパース回帰

中野慎也、樋口知之、地球科学におけるシミュレーションとビッグデータ  
ーデータ同化とエミュレーションー、電子情報通信学会誌、Vol.97(10),  
pp.869-875, 2014.  
樋口知之、中村和幸、データ同化によるオンラインセンシングの高度化、  
計測自動制御学会誌、Vol.51(9), 2012.  
長尾大道、佐藤光三、樋口知之、マルコフ連鎖モンテカルロ法を利用した  
トレーサー試験からフラクチャーの物理パラメータを推定する方法、  
石油技術協会誌、Vol.78(2), pp.197-209, 2013.  
Iba, Y. and Akaho, S., Gaussian process regression with measurement error,  
IEICE Trans. E93-D(10), 2010.

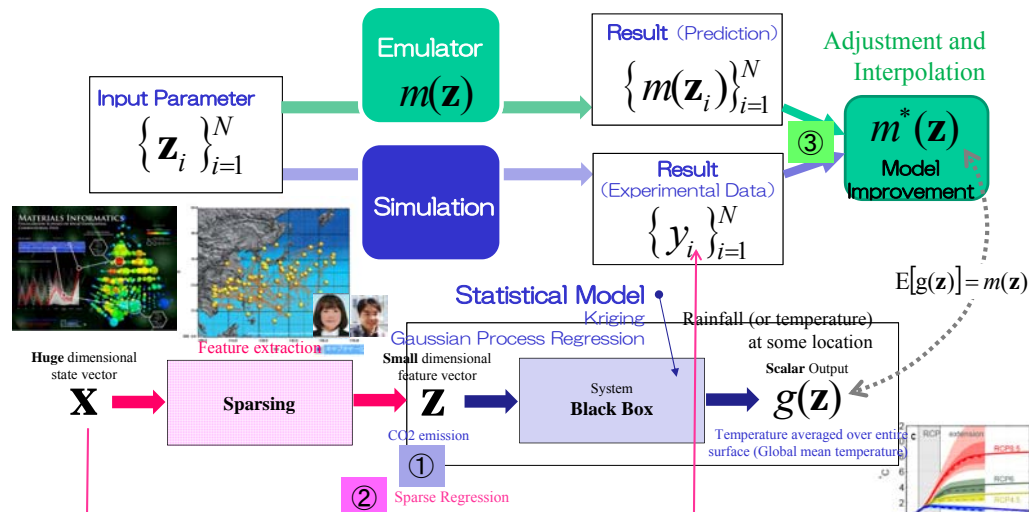


25/42

## 2.b 構成法

26/42

Simulation → Data Assimilation → Emulator : Monte Carlo Experiments  
Statistical model for predicting an output given input parameters → Experimental Design, UQ

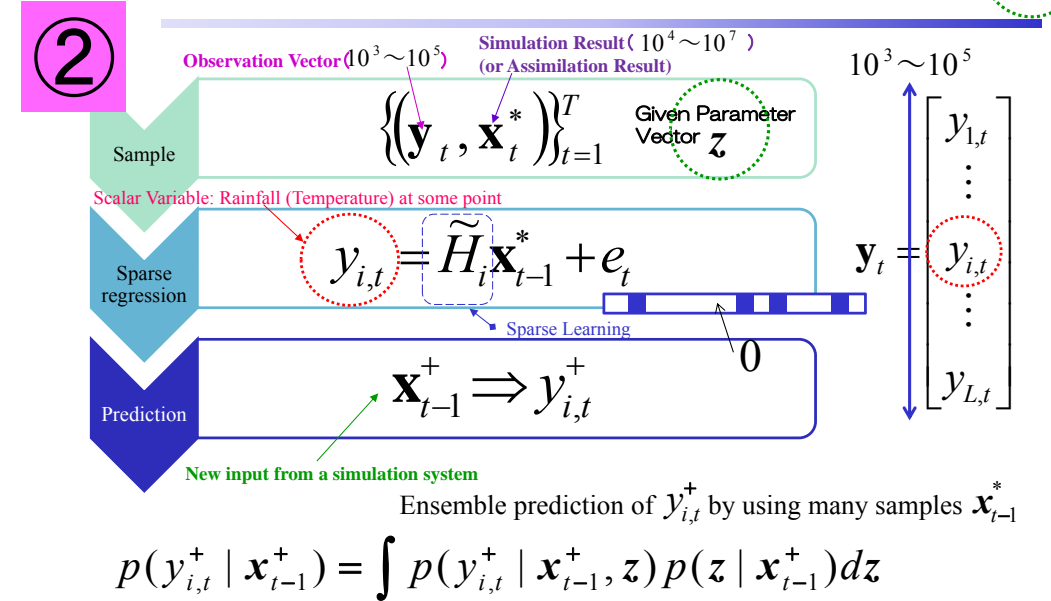
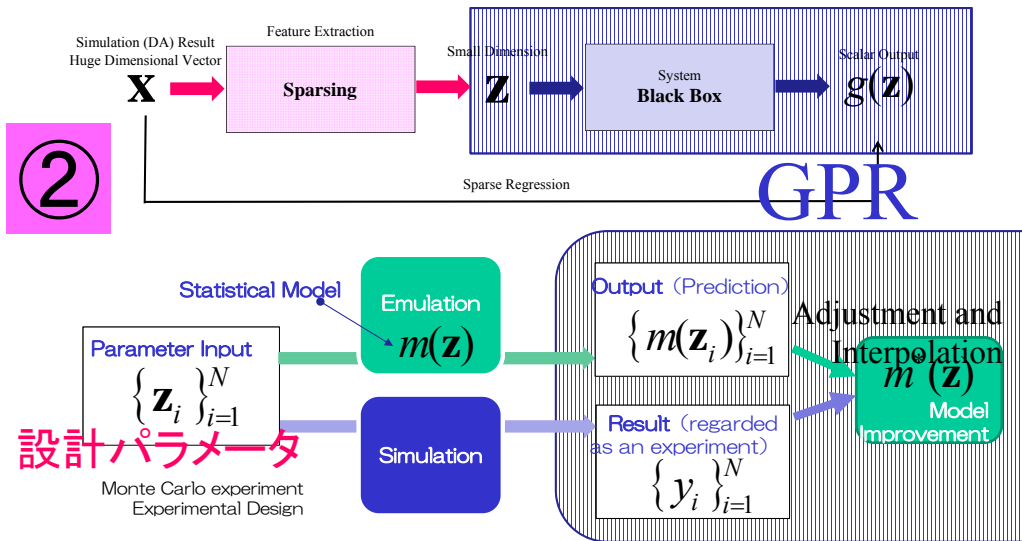


27/42

## エミュレータの設計法

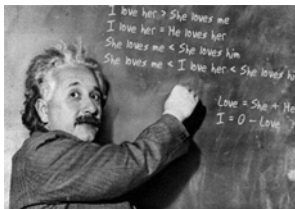
- ① 目的変数(スカラー値)を決める。
- ② シミュレーションを複数回走らせる。  
3. スパース回帰により、入力パラメータベクトルを同定する。
- ③ 4. GPRにより、出力関数(応答局面)をもとめる。

28/42



### 3. スパース回帰

LASSO, CS, NMF



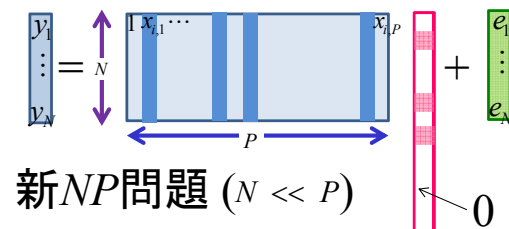
### スパース（疎性を利用した）最適化

例：多変量回帰

$$y_i = a_0 + a_1 x_{i,1} + a_2 x_{i,2} + \dots + a_p x_{i,p} + e_i \quad (i=1, \dots, N)$$

$$\mathbf{y} = \mathbf{H}\mathbf{a} + \mathbf{e}$$

モデルでは表現できない部分  
(誤差というの是不適切)



新NP問題 ( $N \ll P$ )

$$\mathbf{a}^* = \min_{\mathbf{a}} \left\{ \|\mathbf{y} - \mathbf{H}\mathbf{a}\|^2 + \lambda \|\mathbf{a}\| \right\}$$

$$\min_{\mathbf{a}} \|\mathbf{y} - \mathbf{H}\mathbf{a}\|^2 \text{ subj. to } \|\mathbf{a}\| \leq \alpha$$



## 解の一意性と解のロバスト化

初期のベイズ型  
逆問題一般型

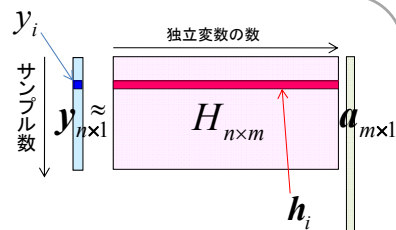
$$\min_a \|y - Ha\|_2^2 + \|Ua\|_{Q/R}^2$$

通常の線形回帰モデル

$$y_i = h_i \cdot a + w_i$$

■  $n > m$  の場合はたいだい大丈夫

■  $m > n$  の場合は解を一意的に定められない。



リッジ回帰  $\min_c \|y - Ha\|_2^2 + \lambda \|a\|_2^2$

ロバストな解を  
求めるために

$$\min_c \|y - Ha\|_2^2 + \lambda \|a\|_1$$

$$\min_c \|y - Ha\|_1 + \lambda \|a\|_1$$

33/42

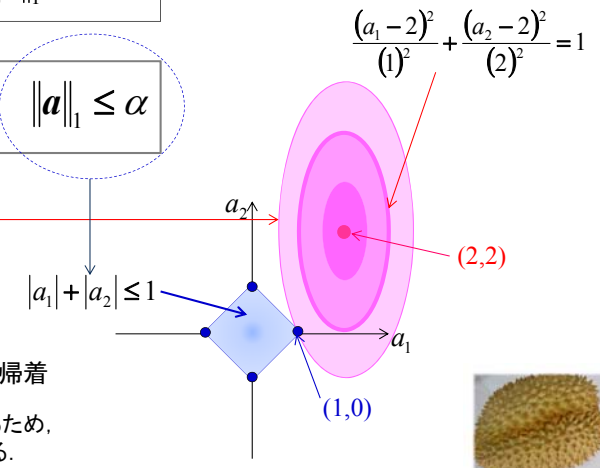
## LASSO: Least Absolute Shrinkage and Selection Operator

■  $m > n$  の場合は解を一意的に定められない。

R. Tibshirani (JRSS B, 1996)

$$\min_a \|y - Ha\|_2^2 + \lambda \|a\|_1$$

$$\min_a \|y - Ha\|_2^2 \quad \text{subj. to} \quad \|a\|_1 \leq \alpha$$



具体的解法は、二次計画問題に帰着

実際には  $\alpha$  をいろいろ変える必要があるため、  
効率的なアルゴリズムが提案されている。

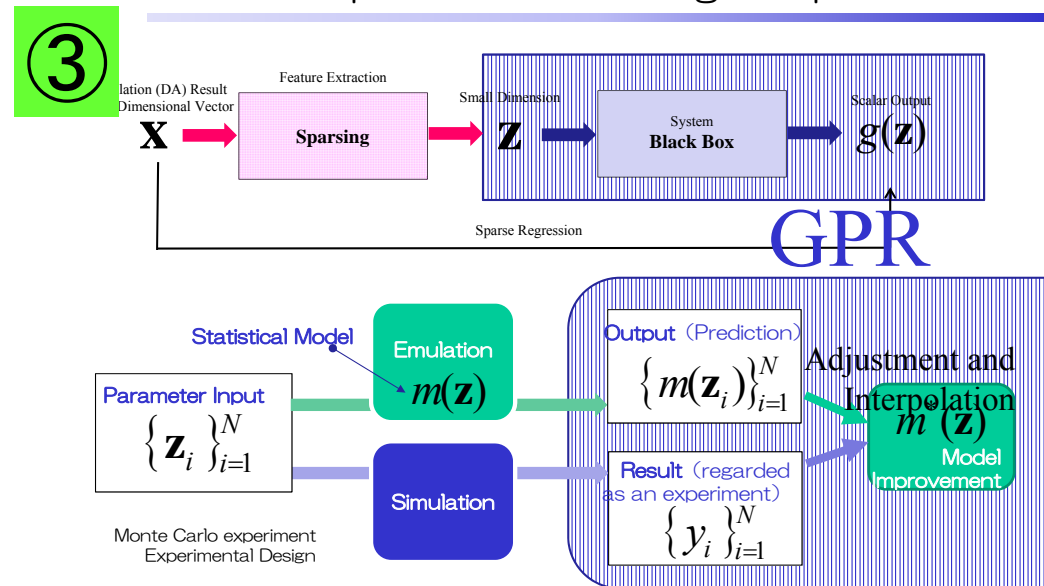
LARS: Least Angle Regression (Efron et al., 2004)

34/42

## 4. Gaussian Process Regression

### Emulator and Emulation

Emulate an output of simulation given parameters



35/42

36/42

## Emulator : Statistical Model (Linear Regression+GP)

### System Model

Gaussian Process:  $g$  is a continuous function of  $\mathbf{z}$

$$p(g) = p(g(\mathbf{z})) = N(m(\mathbf{z}), k(\mathbf{z}, \mathbf{z}') \cdot \tau^2)$$

$$E[g(\mathbf{z})] = m(\mathbf{z}) = H \begin{pmatrix} 1 \\ \mathbf{z} \end{pmatrix} \quad \text{Cov}[g(\mathbf{z}), g(\mathbf{z}')] = k(\mathbf{z}, \mathbf{z}') \cdot \tau^2$$

Regression Coefficient

Statistical model is given

Kernel function

$$k(\mathbf{z}, \mathbf{z}') = \exp\{-(\mathbf{z} - \mathbf{z}')^T R(\mathbf{z} - \mathbf{z}')\}$$

Kernel function: Mutual Distance between input parameters

Common variance

### Observation Model

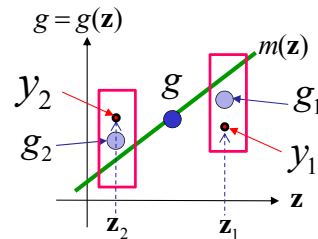
$$y_i = g_i + e, \quad e \sim N(0, \sigma^2) \quad (i=1, \dots, N)$$

Posterior Distribution

$g$  is a continuous function

$$p(g | Y) \propto p(Y | g) p(g)$$

$$Y^T = [y_1, \dots, y_N]$$



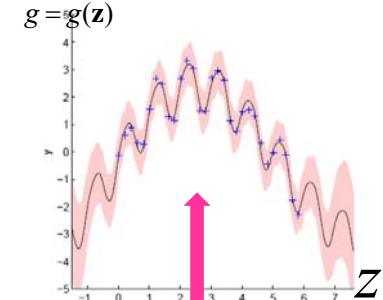
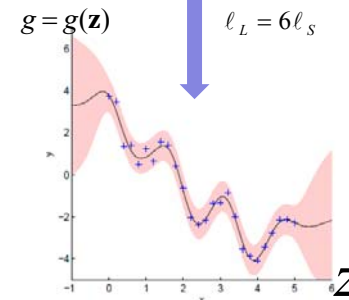
37/42

## Emulator : Design of Kernel function

$$\text{Cov}_y[g(\mathbf{z}), g(\mathbf{z}')] = \tau^2 k(\mathbf{z}, \mathbf{z}') + \sigma^2 \delta(\mathbf{z} - \mathbf{z}') \quad E_y[g(\mathbf{z})] = m(\mathbf{z})$$

$$\tau_1^2 \exp\left[-\frac{(z - z')^2}{2\ell_L^2}\right] + \tau_2^2 \exp\left[-\frac{(z - z')^2}{2\ell_S^2}\right] + \sigma^2 \delta(z - z')$$

Kernel function is NOT related to a mean function  $m(\mathbf{z})$ .

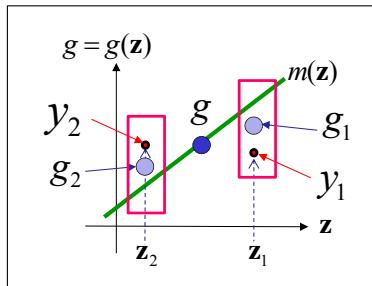


Gaussian Processes for Regression: A Quick Introduction  
M. Ebdon, August 2008

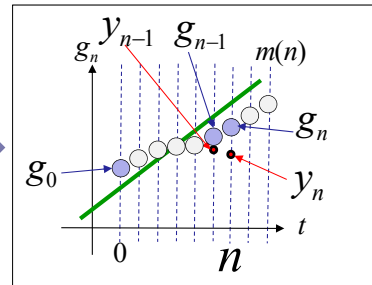
$$\tau_1^2 \exp\left[-\frac{(z - z')^2}{2\ell^2}\right] + \tau_2^2 \exp[-2 \sin^2(\nu\pi(z - z'))] + \sigma^2 \delta(z - z')$$

38/42

## Emulator : Similar structure to SSM if discretizing



If  $g$  is defined only on the discrete point



$$\begin{cases} g_n = m(n) + g_{n-1} + v_n, & v_n \sim N(0, \tilde{\tau}^2) \\ y_n = g_n + e, & e \sim N(0, \sigma^2) \end{cases}$$

Kernel function should positive definite

$$k(n, n-1) = \exp\{-(1)^T R(1)\}$$

39/42

## Emulator : Online Adjustment (Calibration) and Interpolation

or simplicity

a case for  $\frac{\sigma^2}{\tau^2} \rightarrow 0$

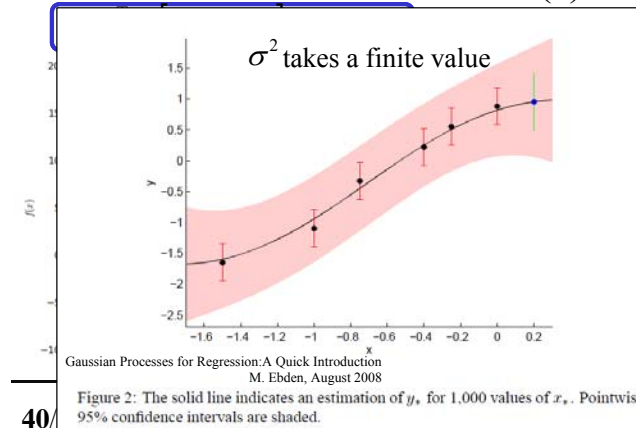
Posterior distribution

$$p(g(\mathbf{z}) | Y) \sim N(m^*(\mathbf{z}), k^*(\mathbf{z}, \mathbf{z}') \tau^2)$$

$$M^T = [m(\mathbf{z}_1), \dots, m(\mathbf{z}_N)]$$

$$m^*(\mathbf{z}) = m(\mathbf{z}) + (Y - M) K^{-1} \mathbf{t}(\mathbf{z})$$

Difference between input  $\mathbf{z}$  and some points



Gaussian Processes for Regression: A Quick Introduction  
M. Ebdon, August 2008

Figure 2: The solid line indicates an estimation of  $y_*$  for 1,000 values of  $x_*$ . Pointwise 95% confidence intervals are shaded.

$$k(\mathbf{z}, \mathbf{z}_1), \dots, k(\mathbf{z}, \mathbf{z}_N)$$

$$k(\mathbf{z}, \mathbf{z}') - \mathbf{t}^T(\mathbf{z}) K^{-1} \mathbf{t}(\mathbf{z}')$$

Distance between Data and Statistical Model

Normalization factor

$$\begin{bmatrix} k(\mathbf{z}_1, \mathbf{z}_1) & \dots & k(\mathbf{z}_1, \mathbf{z}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{z}_N, \mathbf{z}_1) & \dots & k(\mathbf{z}_N, \mathbf{z}_N) \end{bmatrix}$$

Distance from data always reduce uncertainty

40/

For simplicity  
In a case for  $\frac{\sigma^2}{\tau^2} \rightarrow 0$

Posterior distribution

$$p(g(\mathbf{z}) | Y, \tau^2, H) \sim N(m^*(\mathbf{z}), k^*(\mathbf{z}, \mathbf{z}')\tau^2)$$

$$p(g | Y) = \int p(g | Y, \tau^2, H) p(\tau^2) p(H) d\tau^2 dH$$

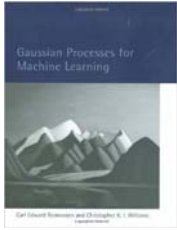


$$p(\tau^2) \propto \frac{1}{\tau^2}, \quad p(H) \propto 1$$

Gamma distribution      Number of samples

$t$  distribution with a degree of freedom  $N - \lambda$   
Emulator  
 $\lambda = \dim(\mathbf{z}) + 1$

For a part of Regression model



#### 参考文献

J. Sacks *et al.*, "Design and analysis of computer experiments," *Statistical Science*, 1989.  
M. C. Kennedy and A. O'hagan, "Bayesian calibration of computer models," *J. Roy. Statist. Soc. Ser. B*, 2001.  
中野、樋口、"地球科学におけるシミュレーションとビッグデータデータ同化とエミュレーション、" 信学会会報、Vol. 97, No.10. 869-875, 2014.  
Rasmussen, C. and C. Williams, *Gaussian Processes for Machine Learning*, MIT Press, 2006.

①

1. 目的変数(スカラー値)を決める。

②

2. シミュレーションを複数回走らせる。

3. スパース回帰により、入力パラメータベクトルを同定する。

③

4. GPRにより、出力関数(応答局面)をもとめる。