

簡便エミュレーションによる 実験計画のスマート化

樋口知之

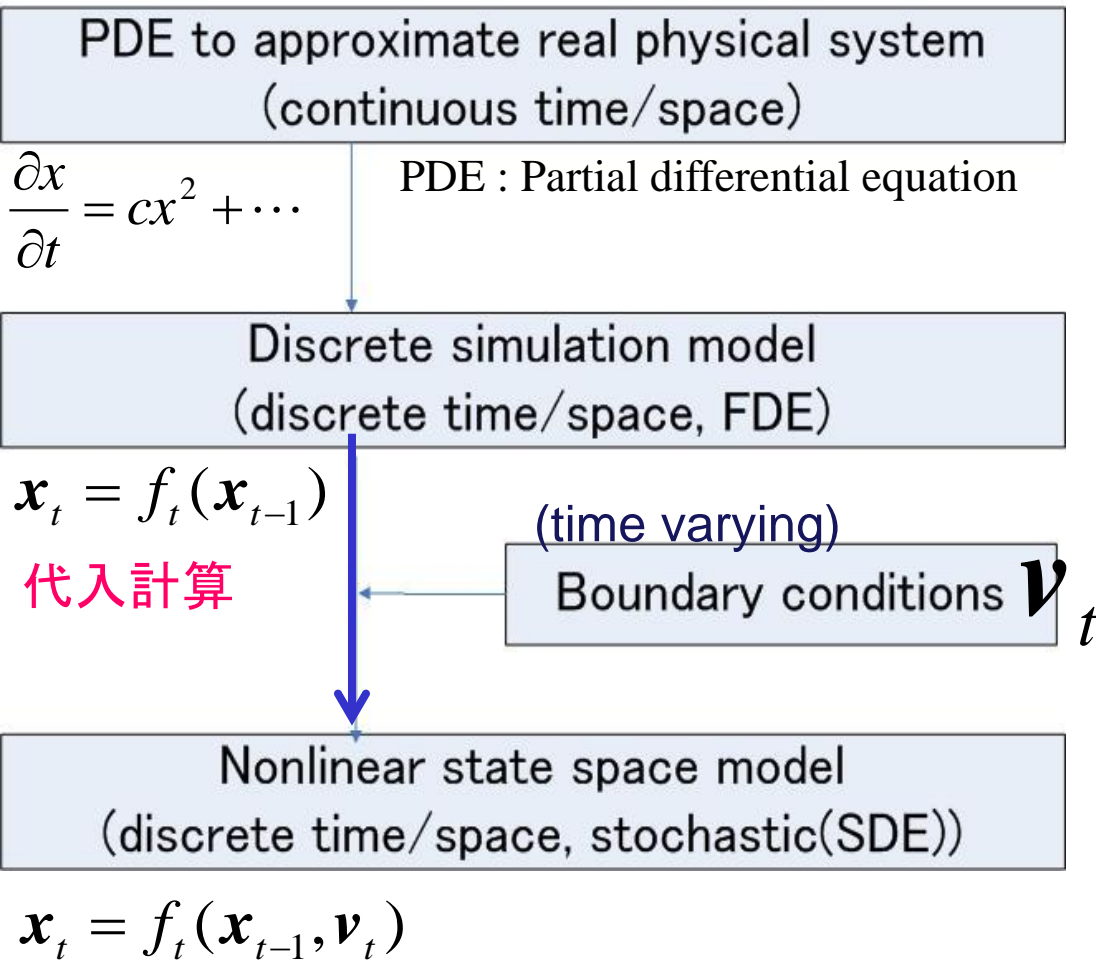
情報・システム研究機構 統計数理研究所

アウトライン

1. データ同化
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
4. GPR: Gaussian Process Regression

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)



State Vector

$$\mathbf{x}_t = \begin{bmatrix} \xi_{1,t} \\ \vdots \\ \xi_{m,t} \\ \xi_{m+1,t} \\ \vdots \\ \xi_{M,t} \\ \theta \end{bmatrix}$$

データ同化と一般状態空間モデル

State Vector (Simulation variables)

システムモデル

Stochastic simulation model

$$\mathbf{x}_t = f_t(\mathbf{x}_{t-1}, \mathbf{v}_t), \quad \mathbf{v}_t \sim p(\mathbf{v} | \boldsymbol{\theta}_{\text{sys}})$$

$$\mathbf{y}_t = h_t(\mathbf{x}_t, \mathbf{w}_t), \quad \mathbf{w}_t \sim p(\mathbf{w} | \boldsymbol{\theta}_{\text{obs}})$$

気象・海洋のデータ
同化の枠組み

$$\mathbf{y}_t = H_t \mathbf{x}_t + \mathbf{w}_t, \quad \mathbf{w}_t \sim N(0, R_{\text{obs}})$$

Observation model

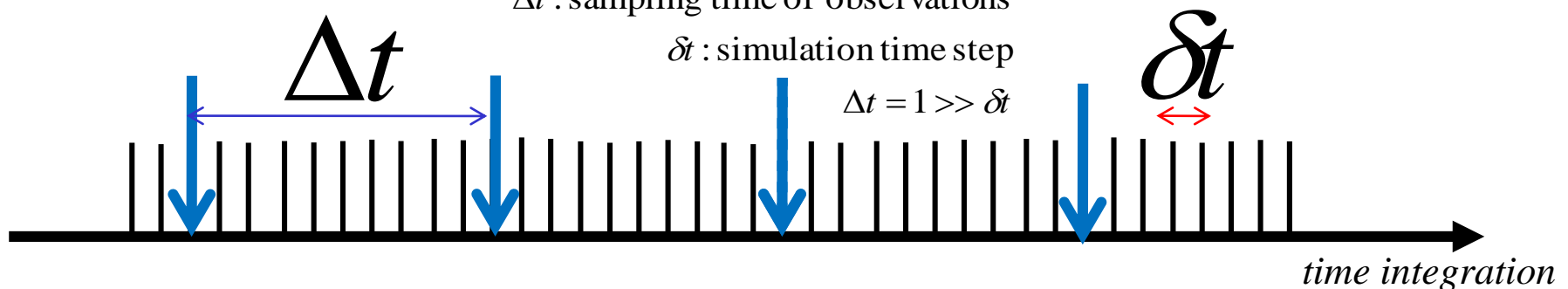
Measurement model

観測モデル

Δt : sampling time of observations

δt : simulation time step

$\Delta t = 1 \gg \delta t$



ベイズの反転公式

x : 興味のある対象

順解析: シミュレーション等

y : データ

ベイズの反転公式

$$p(\underline{x} | \underline{y}) = \frac{p(\underline{y} | \underline{x}) p(\underline{x})}{\sum p(\underline{y} | \underline{x}) p(\underline{x})}$$

逆解析

ベイズの定理。等号の右側と左側で、赤と青で示した変数部分の縦棒との相対関係が反転していることがわかる。この事実により、ベイズの反転公式と呼ばれる。

“連成”ベイズ (ベイズシミュレーション) の普及

- ベイズモデルの階層がますます高層化 (HDP: 階層ディリクレプロセス)
- ノンパラメトリックベイズ法の浸透
 - ・ディリクレプロセスの利用が自然言語の分野ですすんだ。
 - ・モデル選択からモデル混合へ。
- ミクロからマクロへ生成モデルの連続結合 (“連成”) による、いわばフォワード推論がベイズ計算の主流

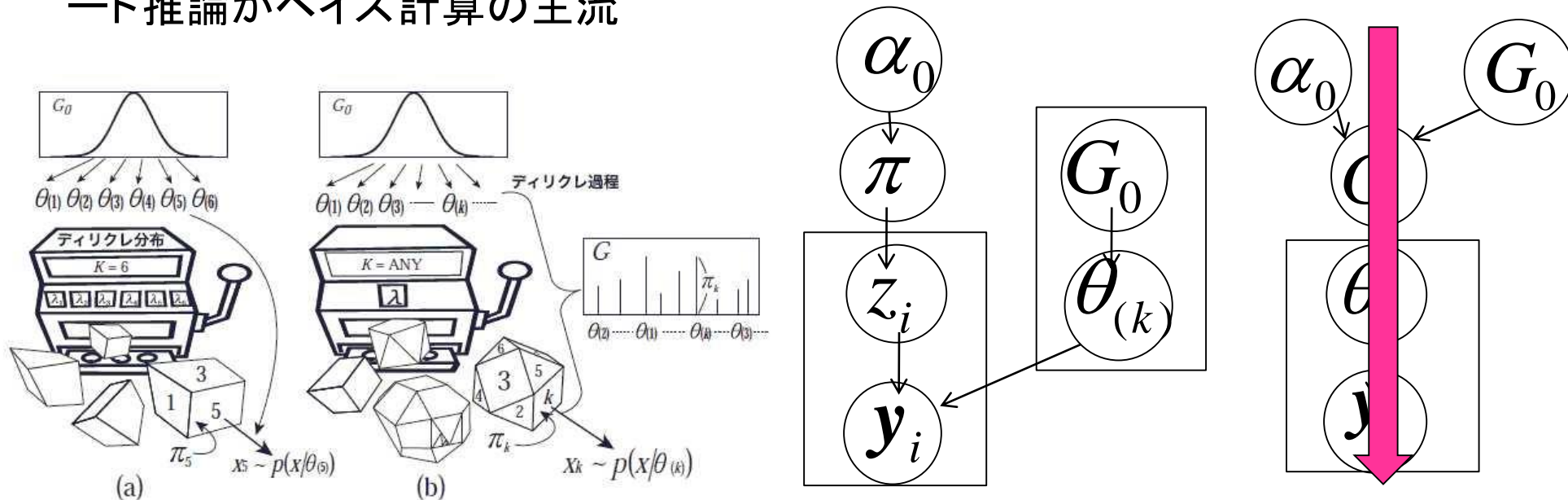


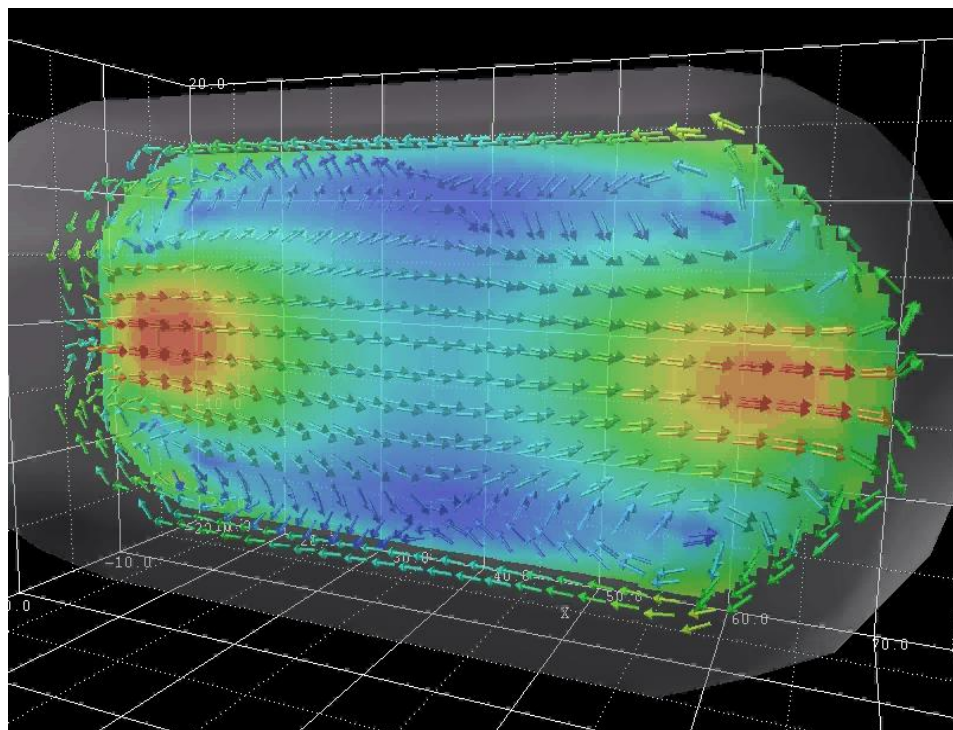
Fig. 2. ディリクレ分布/ディリクレ過程サイコロ生成器

Data Assimilation on Intracellular Fluid in C.elegans embryo

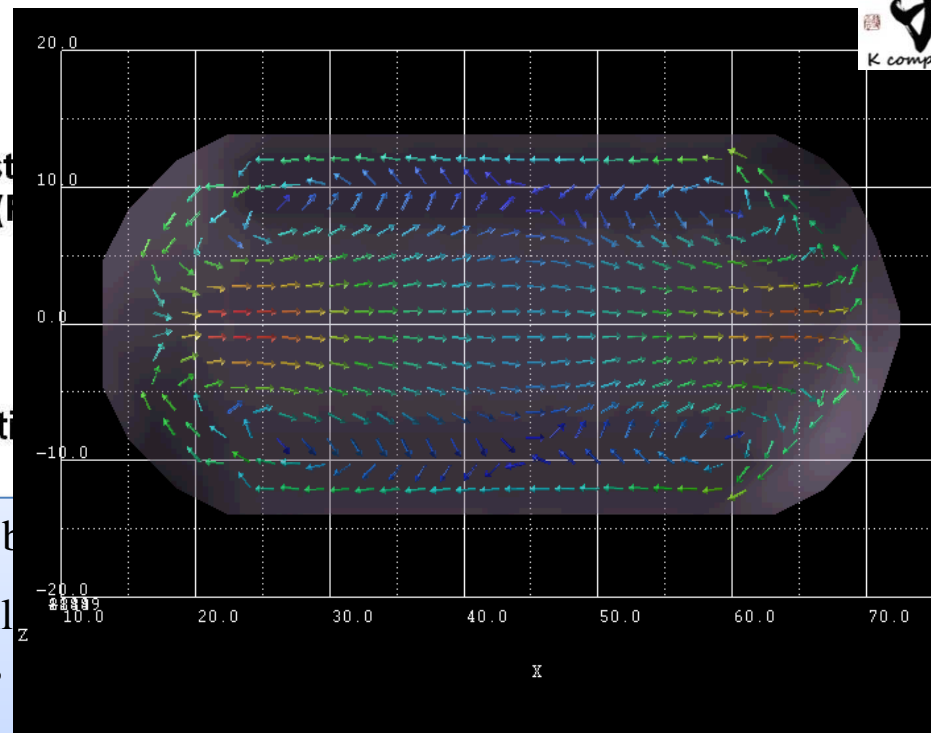
(Collaboration with Dr. Kimuta at NIG)



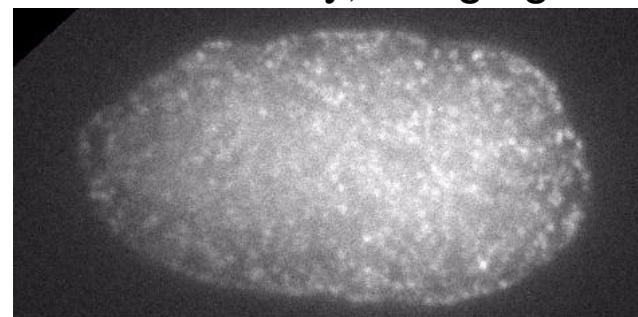
Assimilation



Observation



We assume that proteins “**myosin**” in the cell wall drive physically such intra-cellular flow, and we perform data assimilation by combining a numerical simulation that solves simultaneously these physical equations and the observation obtained by the image analyses.



Cell tracking

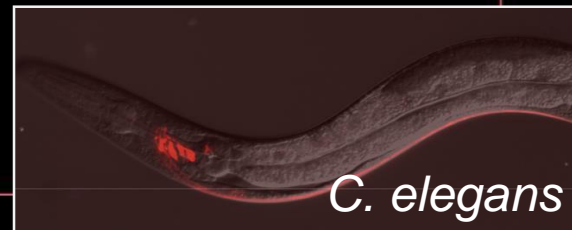


Tracking result

Centroid Cell Positions at time=20

4D imaging system

- The neuronal nuclei of adult *C. elegans* are labeled with a red fluorescent protein sensor called mCherry.
- A combination of the fluorescent protein sensor for calcium ion visualizes the neuronal activity of the multiple neurons in the living state.
- The temporal changes of the neuronal activities are measured by a confocal microscope

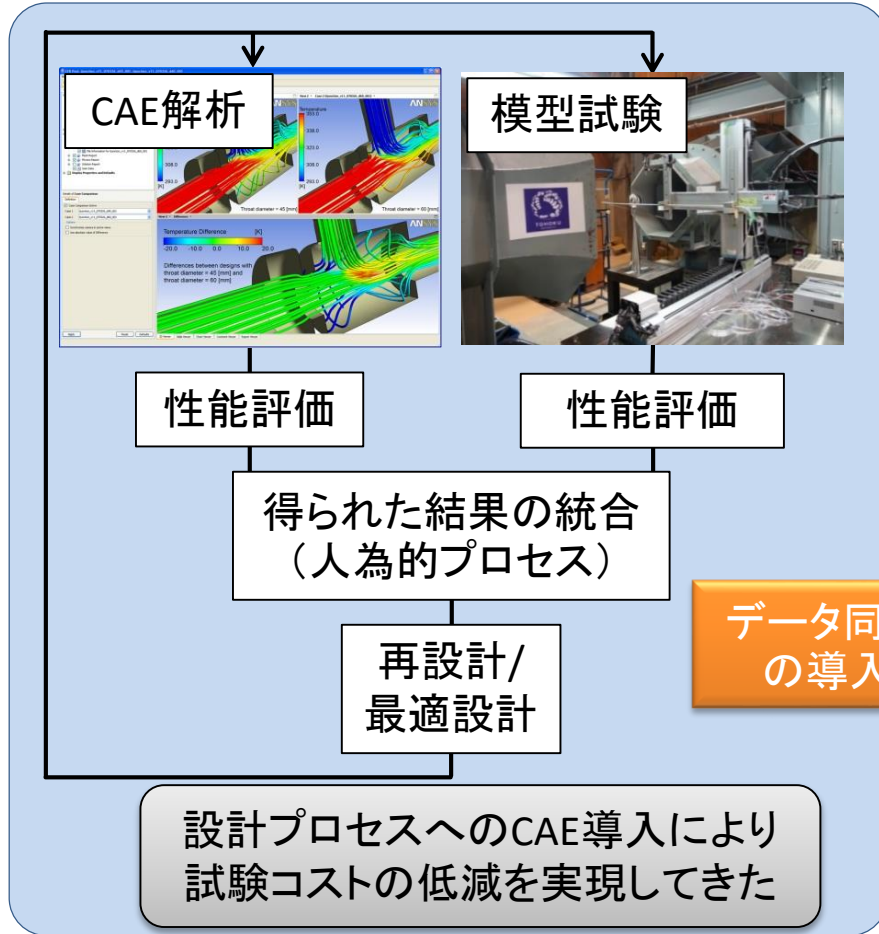


C. elegans

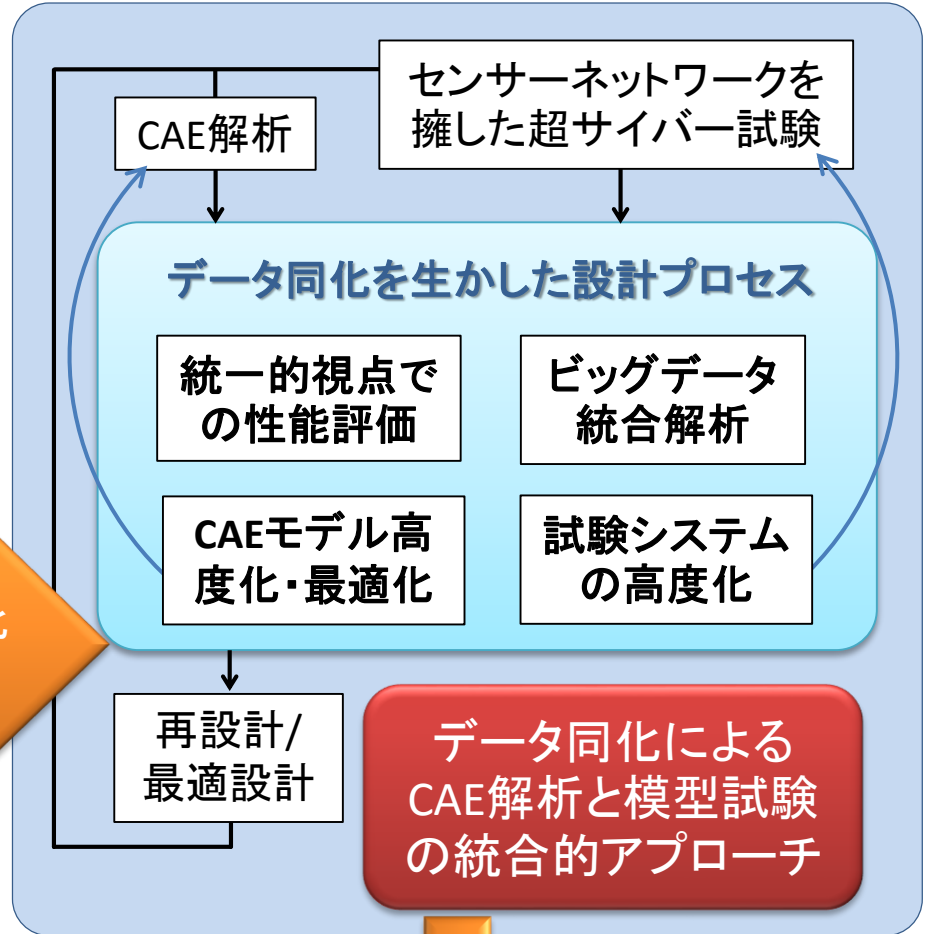


Quantitative analysis of neural activities with Ca²⁺ imaging of *C. elegans* whole brain

データ同化による設計技術の革新



データ同化
の導入



データ同化による設計プロセスで利用される試験・CAEデータの信頼性・精度向上

- 時間コストの低減
- 製品の高度化による競争力強化

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SAMSI(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:

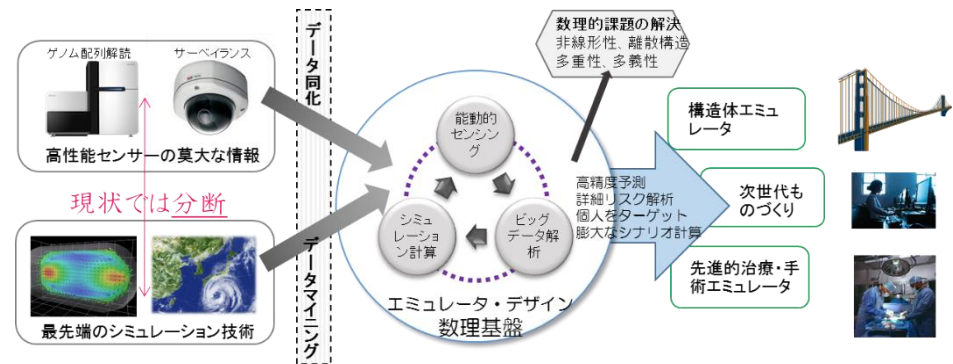
ガウス過程回帰や、その古典版とも言えるクリギング
次元削減を目的としたスパース回帰

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
—データ同化とエミュレーション—、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.



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

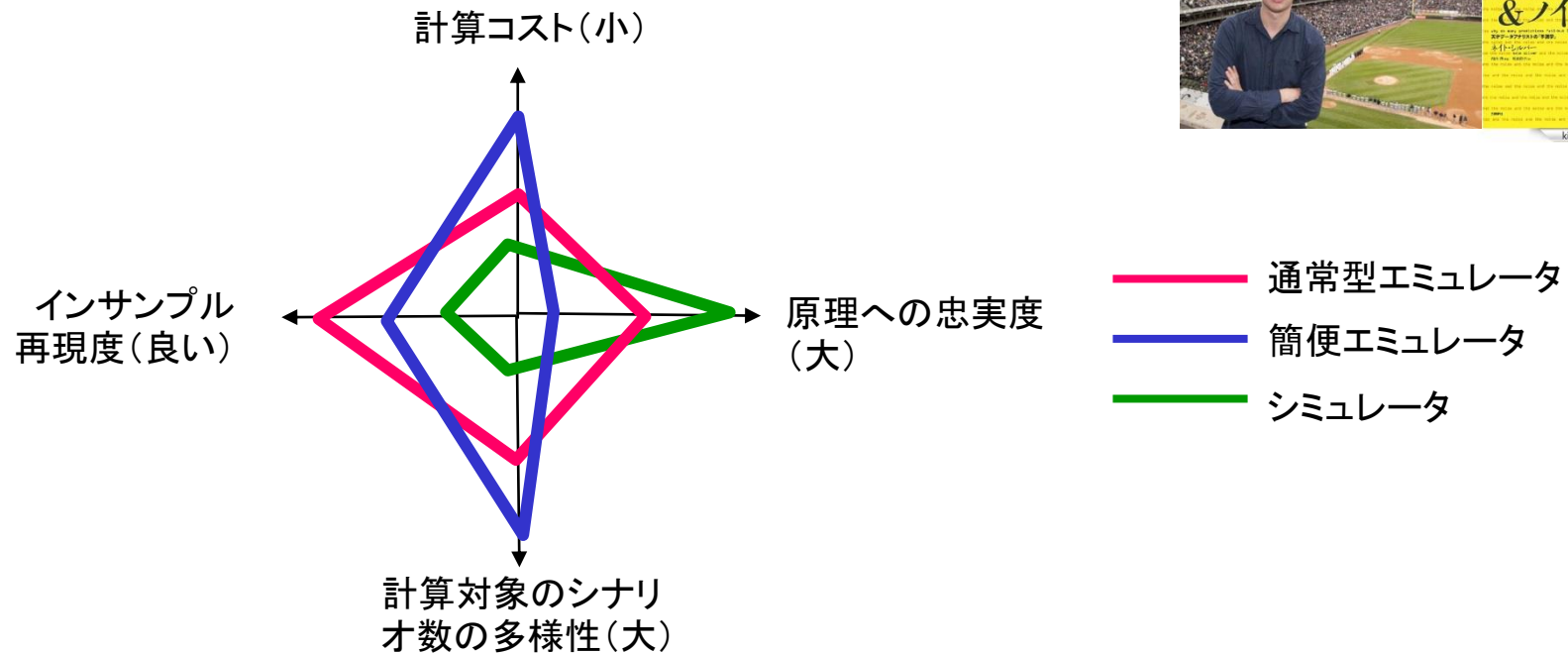
エミュレーションやエミュレータは、模倣対象のシステムにおいて、**予測できる現象より予測できない現象が支配的**である場合に使われる。また、非常に高い安全性が要求される場合にも良く使われる。**予測できる現象が支配的な場合や、完全に模倣することが難しい場合はシミュレーション技術**を使う。

(Wikipediaより)

計算コストと予測精度

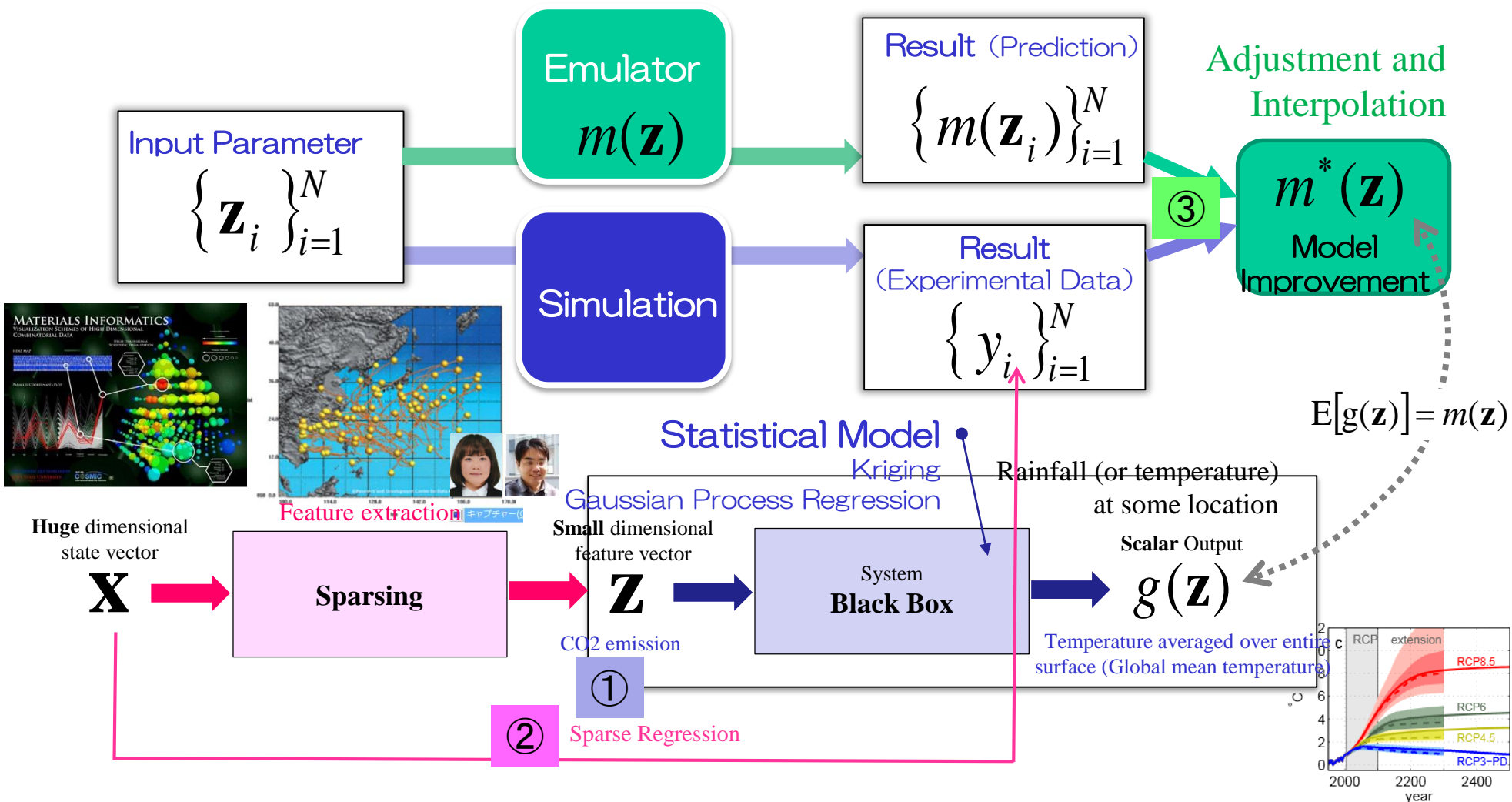
- 高い計算コスト(オンライン計算が難しい)
- 不確実性の取り扱い方

“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。”



Simulation → Data Assimilation → Emulator : Monte Carlo Experiments → Experimental Design, UQ

Statistical model for predicting an output given input parameters



エミュレータの設計法

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。

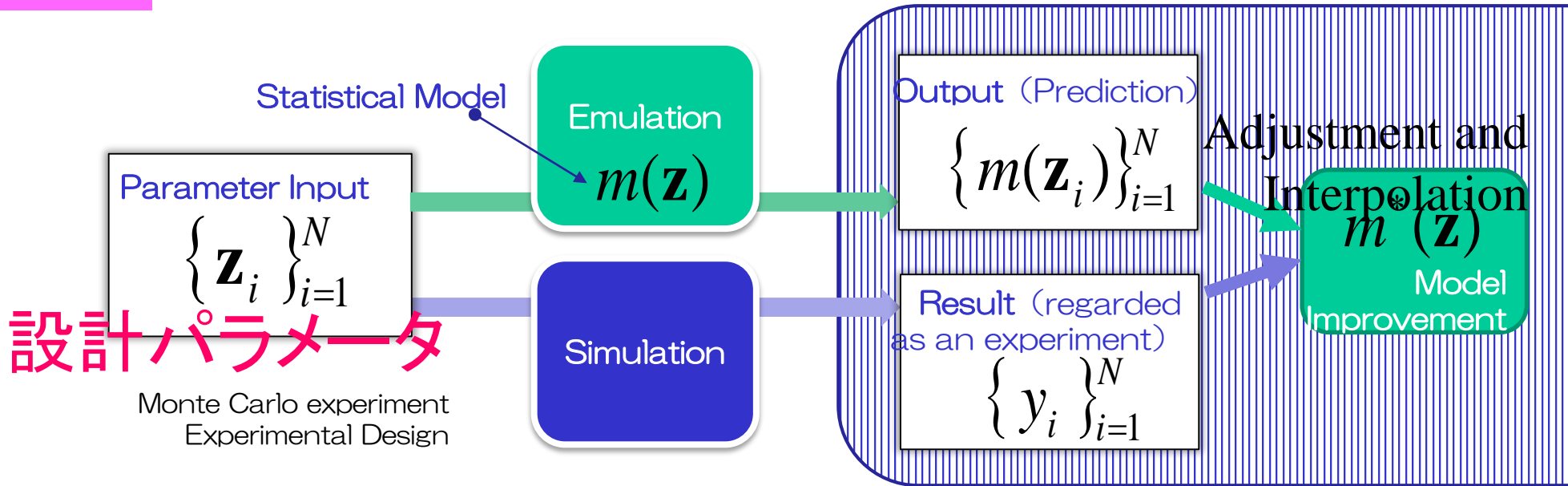
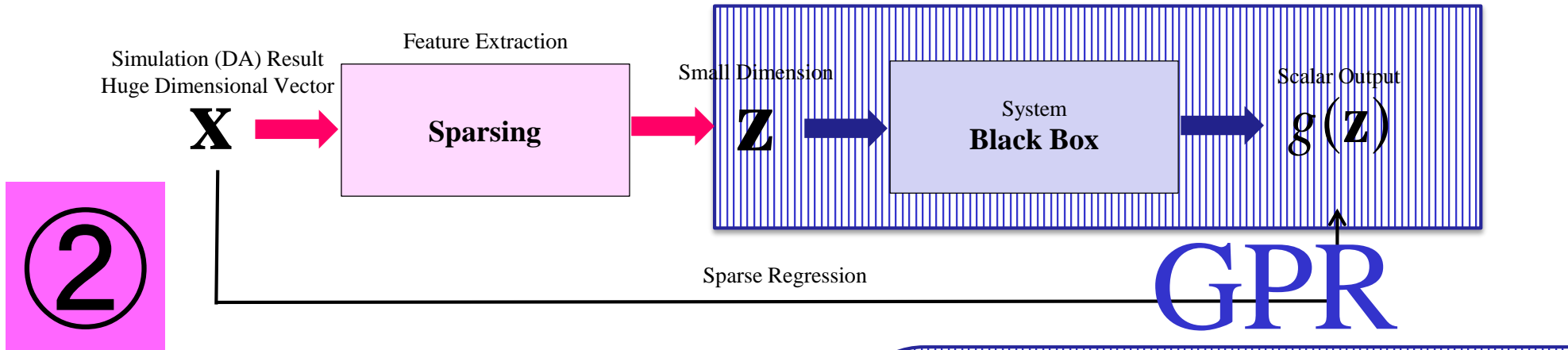
3. スパース回帰により、入力パラメータベクトルを同定する。

③

4. GPRにより、出力関数(応答局面)をもとめる。

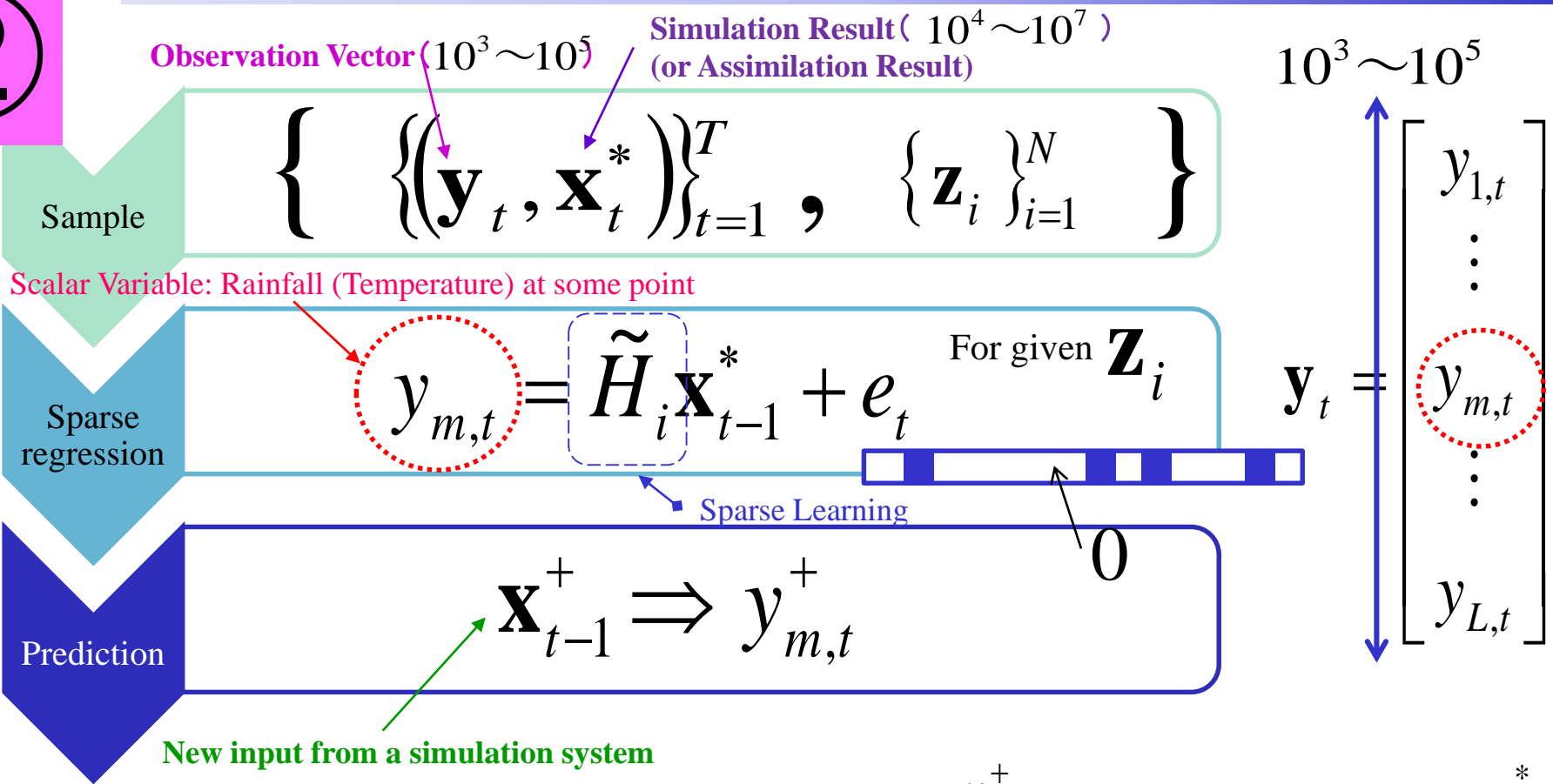
Emulator and Emulation

Emulate an output of simulation given parameters



Climate Prediction : Climate Emulator $p(y_{m,t} \mid \mathbf{x}_{t-1}^+, \mathbf{z})$

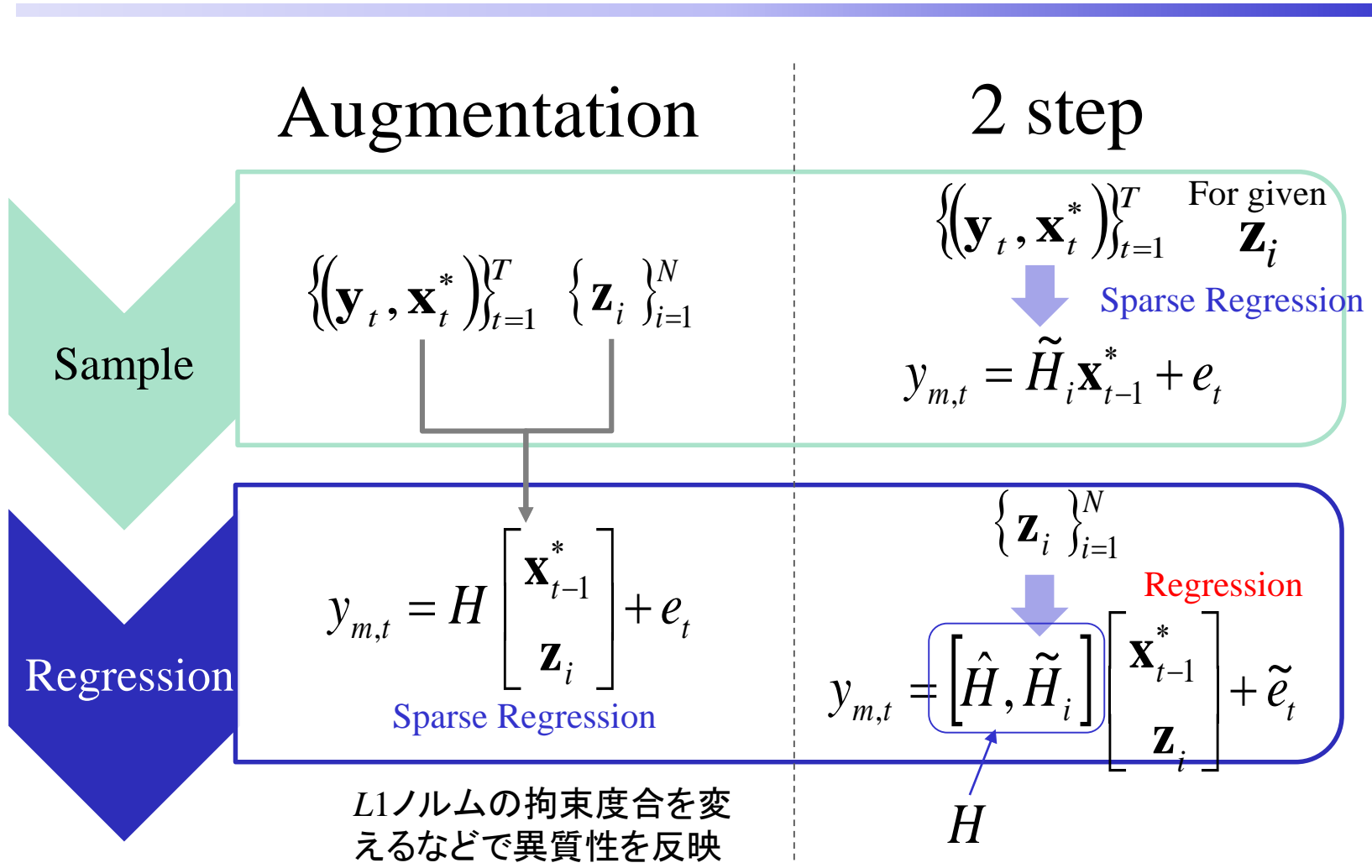
2



Ensemble prediction of $y_{m,t}^+$ by using many samples \mathbf{x}_{t-1}^*

$$p(y_{m,t}^+ \mid \mathbf{x}_{t-1}^+) = \int p(y_{m,t}^+ \mid \mathbf{x}_{t-1}^+, \mathbf{z}) p(\mathbf{z} \mid \mathbf{x}_{t-1}^+) dz$$

②



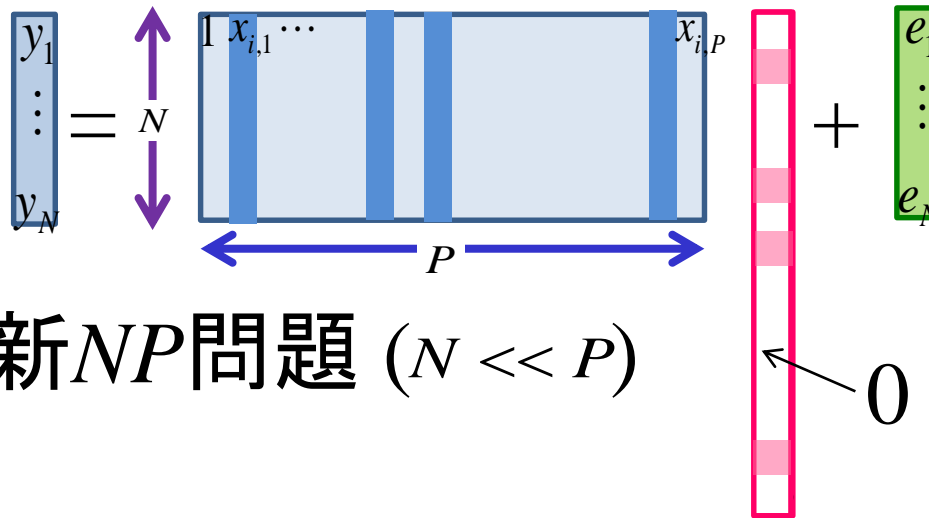
スパース（疎性を利用した）最適化

例：多変量回帰

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

$$y = Ha + e$$

モデルでは表現できない部分
(誤差というのは不適切)



新NP問題 ($N \ll P$)

$$a^* = \min_a \left\{ |y - Ha|^2 + \lambda |a| \right\}$$

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

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

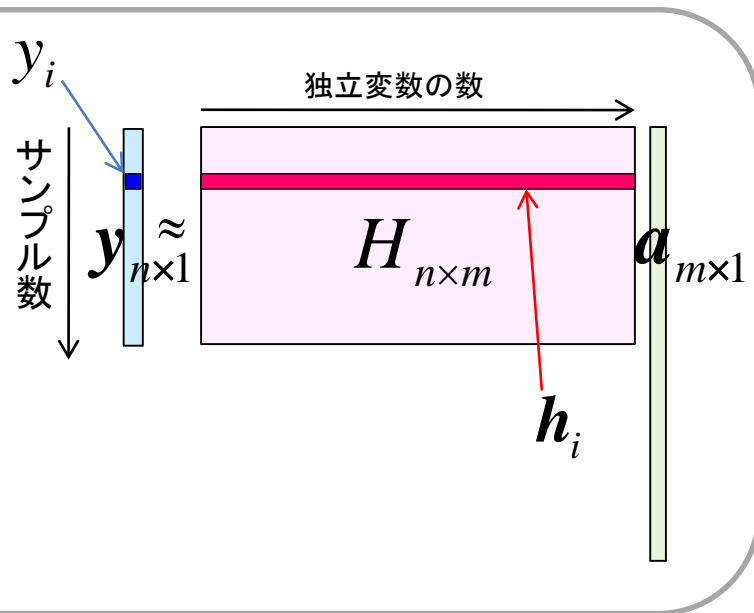
$$\min_a \cdot \|y - Ha\|_2^2 + \|Ua\|_{Q/R}^2$$

通常の線形回帰モデル

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

■ $n > m$ の場合はたいてい大丈夫

■ $m > n$ の場合は解を一意に定められない。



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

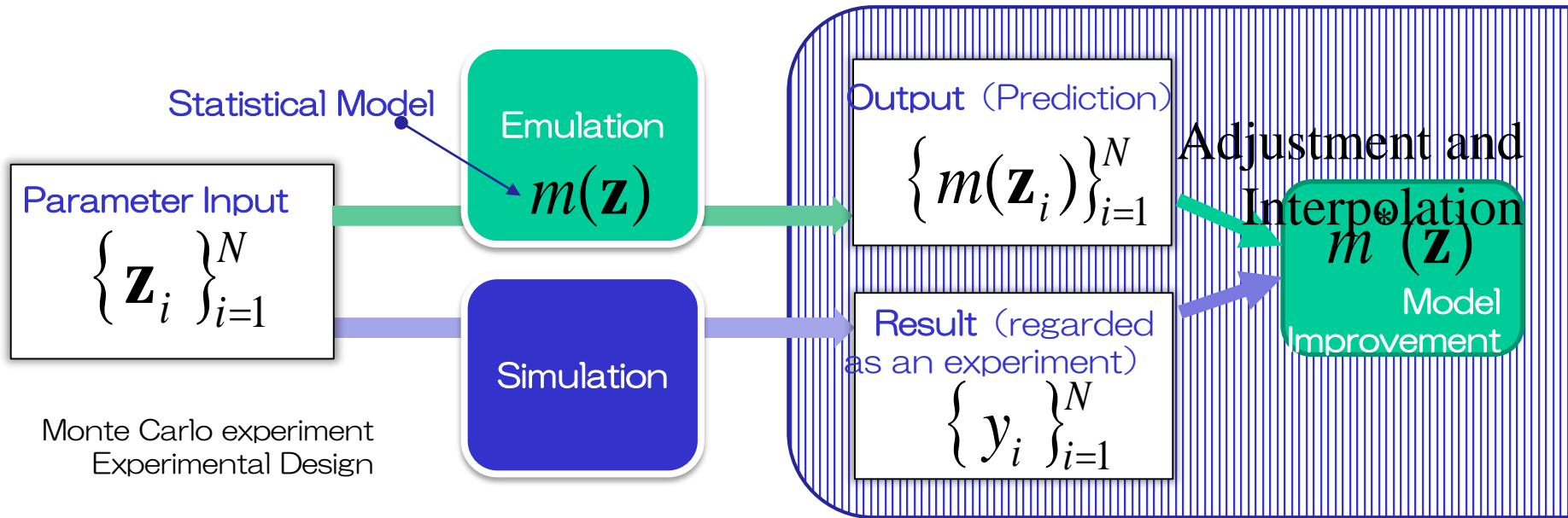
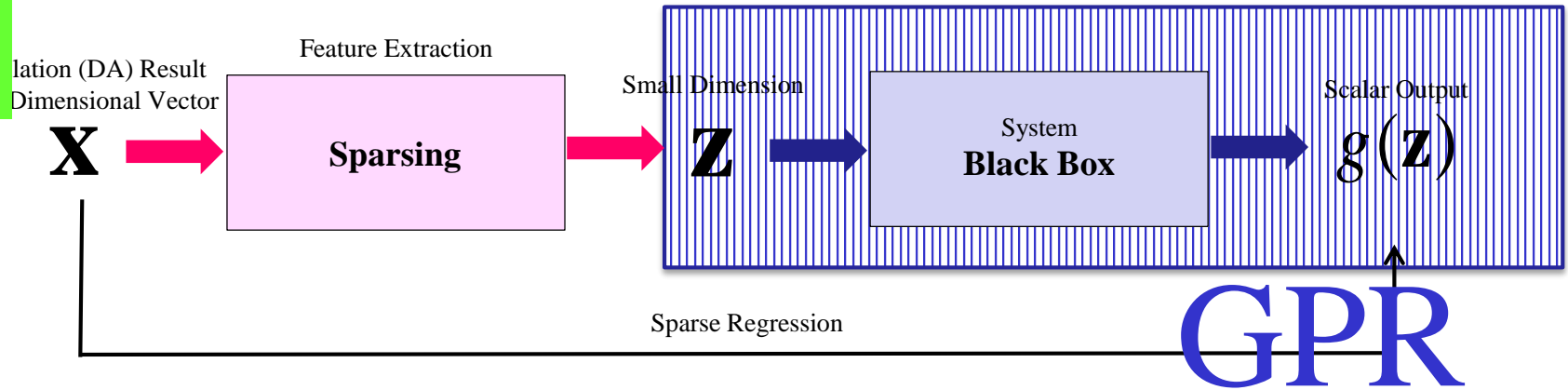
ロバストな解を
求めるために

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

Emulator and Emulation

Emulate an output of simulation given parameters

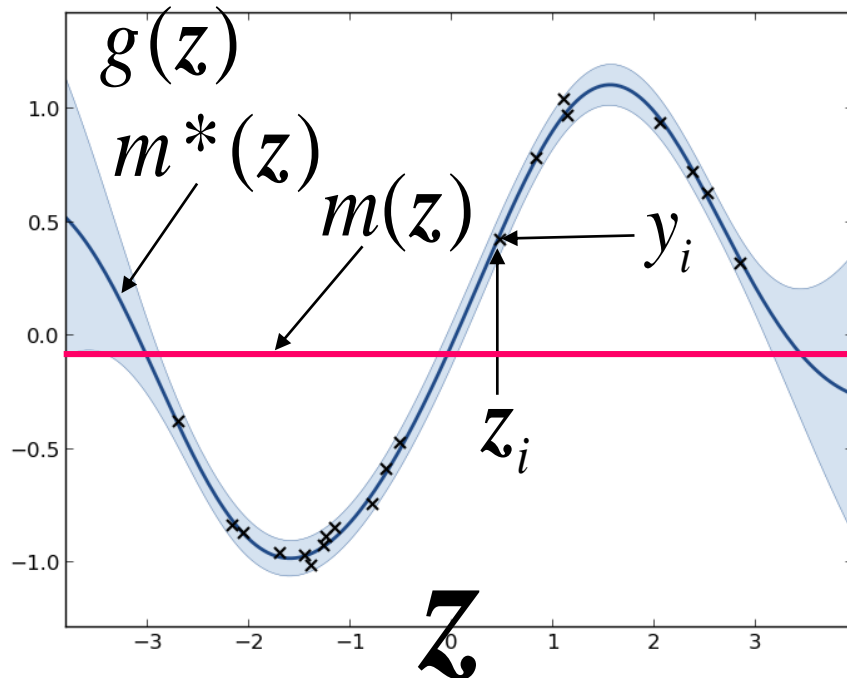
③



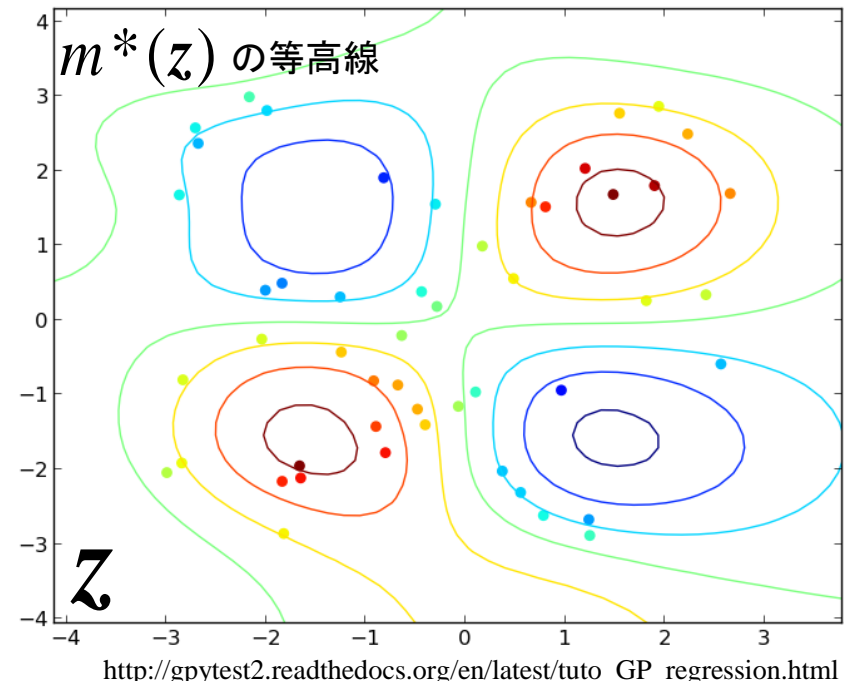
要は内挿の一手法

データ $\{y_i\}_{i=1}^N$: シミュレーション結果
入力 z : 設計パラメータ (多次元)
出力 $g(z)$: 目的変数を出力とする確率的応答関数

1次元



2次元



Emulator : Statistical Model (Linear Regression+GP)

System Model

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

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

$$E[g(\mathbf{z})] = m(\mathbf{z}) = \underset{\text{Regression Coefficient}}{H} \cdot \begin{pmatrix} \mathbf{x}_{t-1}^* \\ \mathbf{z} \end{pmatrix} \quad \text{Cov}[g(\mathbf{z}), g(\mathbf{z}')] = \underset{\text{Covariance function}}{k(\mathbf{z}, \mathbf{z}')} \cdot \underset{\text{Common variance}}{\tau^2}$$

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

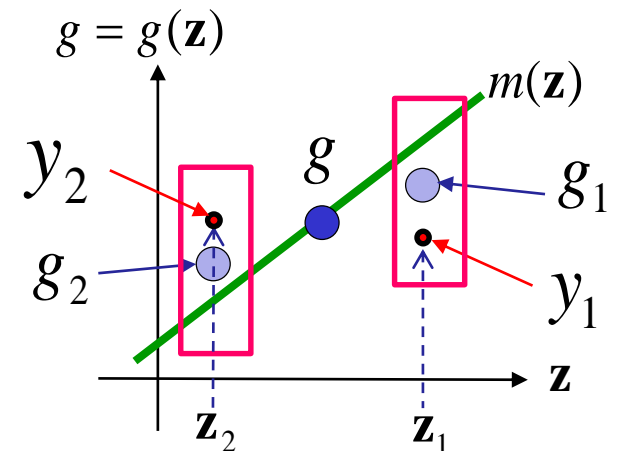
$$y_i = g_i + e, \quad e \sim N(0, \sigma^2) \quad (i=1, \dots, N)$$

Posterior Distribution

g is a continuous function

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

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



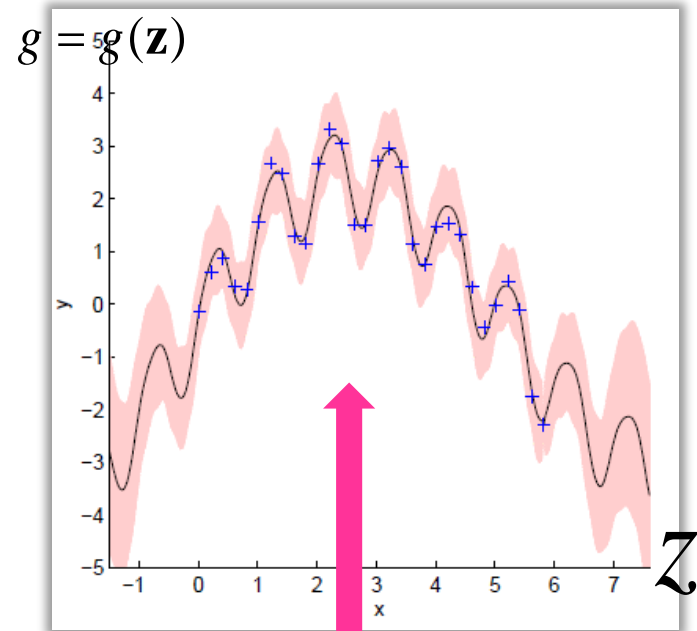
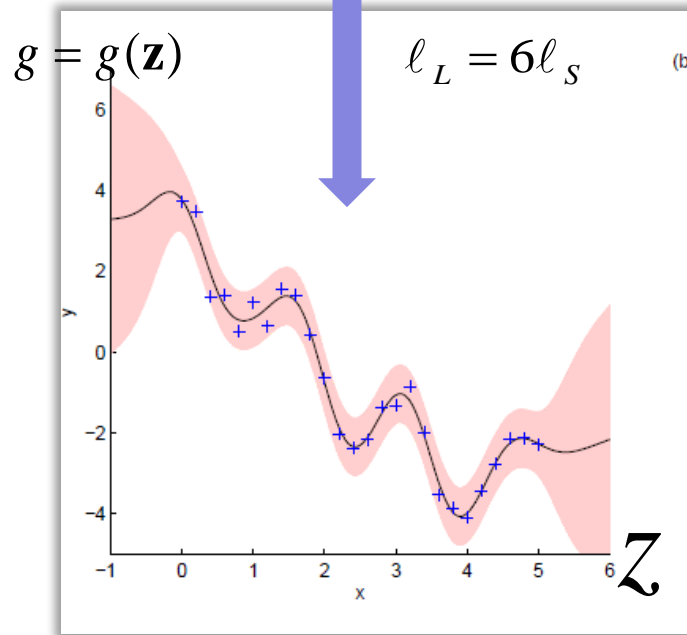
Emulator : Design of Kernel function

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

$$E_y[g(\mathbf{z})] = m(\mathbf{z})$$

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

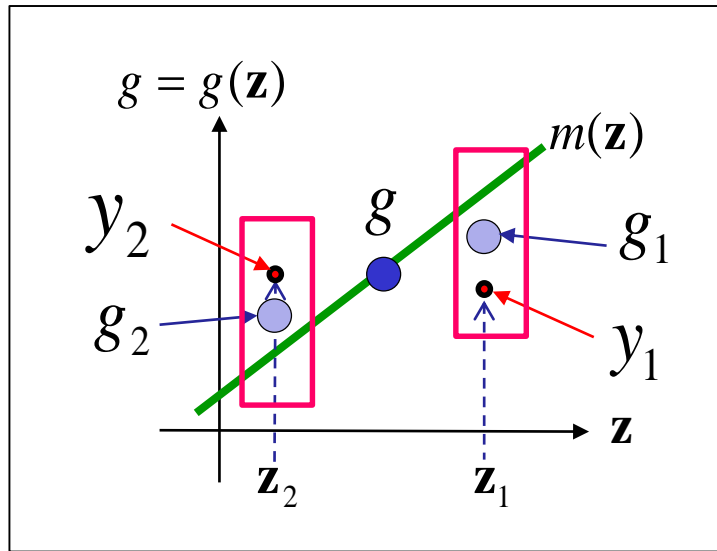
Kernel function is NOT related to a mean function $m(\mathbf{z})$.



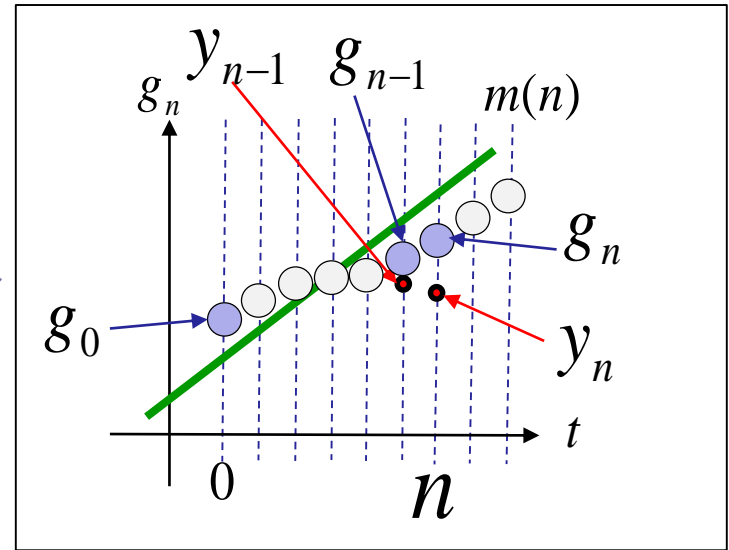
Gaussian Processes for Regression: A Quick Introduction
M. Ebden, August 2008

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

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

Emulator : Online Adjustment (Calibration) and Interpolation

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

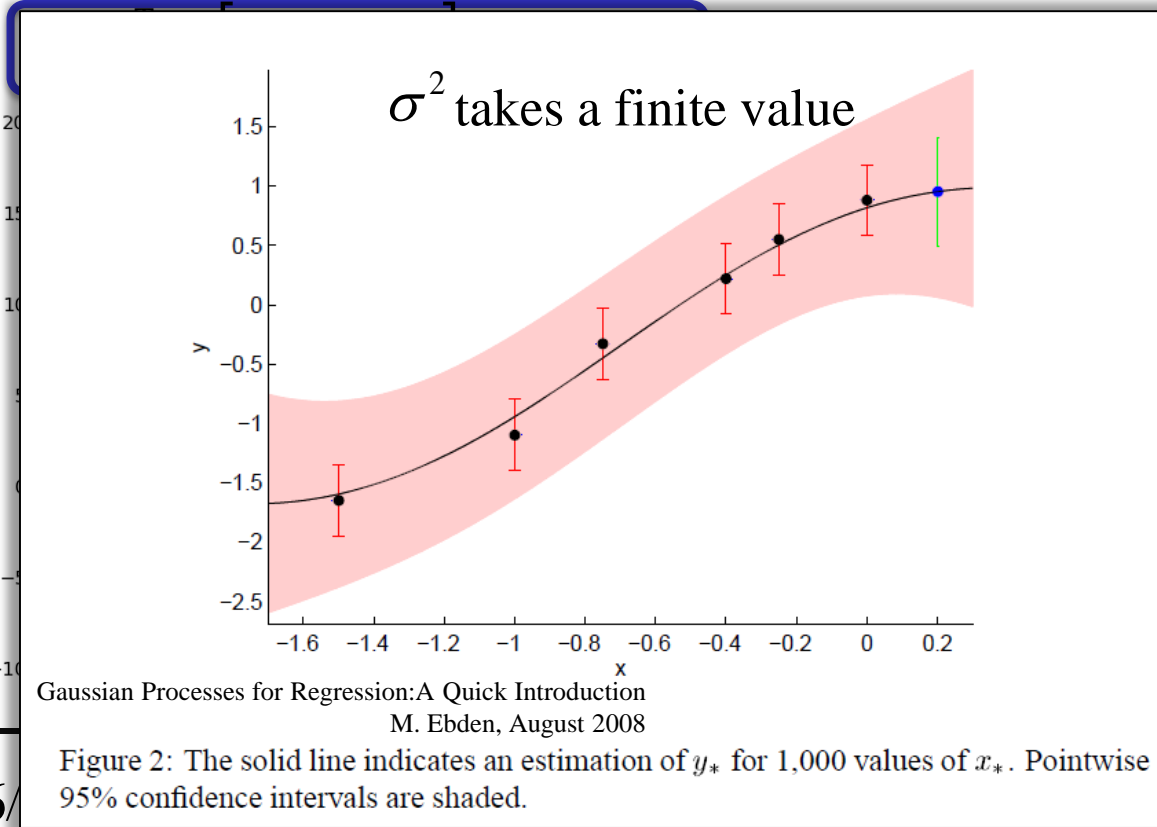
Posterior distribution

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

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

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

Difference between input \mathbf{z} and some points



Difference between Data and Statistical Model

Normalization factor

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

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

Information from data always reduce uncertainty

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

Emulator : Obtain a full Bayes model

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

Posterior distribution

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

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

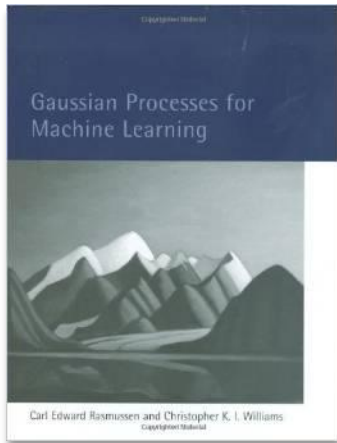


$$p(\tau^2) \propto \frac{1}{\tau^2}, \quad p(H) \propto 1$$

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.

エミュレータの設計法(再掲)

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。