# Predictive Distribution for Gaussian Process Models with Design-Based Subagging

#### Linglin He<sup>a</sup>, Yufan Liu<sup>b</sup> and Ying Hung<sup>a</sup>

<sup>a</sup>Rutgers University, Department of Statistics and Biostatistics <sup>b</sup> Dun & Bradstreet

Oct 13, 2017

1/23

#### Overview

#### Introduction

- 2 LHD-Based Block Bootstrap
- 3 Consistency of the Bootstrap Estimators
- 4 Bootstrap Predictive Distribution
- **5** Numerical Studies
- 6 Future Work
  - **7** Summary

#### **Computer Experiments**

- Computer experiments refer to the study of real systems using complex mathematical models
- They have been widely used as alternatives to physical experiments, which sometimes are unethical, impossible, inconvenient or too expensive.
- They are nearly deterministic in the sense that a particular input will produce almost the same output.
- Therefore, it is desirable to build an interpolator for computer experiment outputs and use it as an emulator for the actual computer experiment.

#### Introduction

#### Gaussian Process Model for Computer Experiments

 Consider a computer experiment which has n inputs x ∈ R<sup>d</sup> and produces univariate output y(x). To analyze the experiments, y(x) is assumed to be a realization from a stochastic process model:

$$Y(\mathbf{x}) = \mu(\mathbf{x}) + Z(\mathbf{x}),$$

- mean function:  $\mu(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$
- $Z(\mathbf{x})$  : stationary Gaussian process with mean 0 and covariance function  $\sigma^2\psi$
- o covariance function:

$$Cov\{Z(\boldsymbol{x}_i), Z(\boldsymbol{x}_j)\} = \sigma^2 \psi(\boldsymbol{x}_i - \boldsymbol{x}_j; \boldsymbol{\theta})$$

where  $\boldsymbol{\theta}$  is a vector of correlation parameter for the correlation function.

#### Gaussian Process Model (Estimation)

• Given *n* observed realizations  $X_n$  and  $Y_n$ , the log-likelihood function, ignoring a constant, can be written as

$$\ell(\boldsymbol{X}_n, \boldsymbol{y}_n, \boldsymbol{\phi}) = -\frac{1}{2\sigma^2} (\boldsymbol{y}_n - \boldsymbol{X}_n \boldsymbol{\beta})^T R_n^{-1}(\boldsymbol{\theta}) (\boldsymbol{y}_n - \boldsymbol{X}_n \boldsymbol{\beta}) - \frac{1}{2} \log |R_n(\boldsymbol{\theta})| - \frac{n}{2} \log(\sigma^2),$$
  
where  $R_n(\boldsymbol{\theta}) = [\psi(\boldsymbol{x}_i - \boldsymbol{x}_j); \boldsymbol{\theta}), i, j = 1, \dots, n]$  is an  $n \times n$  correlation matrix.

• The MLE can be obtained by

$$\hat{\boldsymbol{\beta}}_n = (\boldsymbol{X}_n^T \boldsymbol{R}_n^{-1}(\boldsymbol{\theta}) \boldsymbol{X}_n)^{-1} \boldsymbol{X}^T \boldsymbol{R}_n^{-1}(\boldsymbol{\theta}) \boldsymbol{y}_n,$$

$$\hat{\sigma}_n^2 = (\boldsymbol{y}_n - \boldsymbol{X}_n \hat{\boldsymbol{\beta}}_n)^T R_n^{-1}(\boldsymbol{\theta}) (\boldsymbol{y}_n - \boldsymbol{X}_n \hat{\boldsymbol{\beta}}_n) / n,$$

$$\hat{\boldsymbol{\theta}}_n = \operatorname*{arg\,min}_{\boldsymbol{\theta}} \{ n \log(\hat{\sigma}_n^2) + \log |R_n(\boldsymbol{\theta})| \}.$$

イロト 不得下 イヨト イヨト 二日

#### Gaussian Process Model (Prediction)

Based on the MLEs, we are interested in predicting y<sub>n+1</sub> at an untried new input x<sub>n+1</sub> and quantifying the uncertainty. To achieve this, the conventional plug-in method predicts y<sub>n+1</sub> by a distribution g(x<sub>n+1</sub> | X<sub>n</sub>, Y<sub>n</sub>, φ̂<sub>n</sub>) which is normally distributed with mean

$$\mu(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n, \hat{\phi}_n) = \boldsymbol{x}_{n+1}^T \hat{\boldsymbol{\beta}}_n + \gamma_n(\hat{\boldsymbol{\theta}}_n)^T \boldsymbol{R}_n^{-1}(\hat{\boldsymbol{\theta}}_n)(\boldsymbol{y}_n - \boldsymbol{X}_n \hat{\boldsymbol{\beta}}_n)$$

and variance

$$\sigma^{2}(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_{n}, \boldsymbol{y}_{n}, \hat{\phi}_{n}) = \hat{\sigma}_{n}^{2} \{1 - \gamma_{n}(\hat{\boldsymbol{\theta}}_{n})^{T} R_{n}^{-1}(\hat{\boldsymbol{\theta}}_{n}) \gamma_{n}(\hat{\boldsymbol{\theta}}_{n})\}$$

where  $\gamma_n(\hat{\theta}_n)$  is the correlation between the new observation and the existing data, i.e.  $\gamma_n(\hat{\theta}_n) = [\psi(\mathbf{x}_i - \mathbf{x}_{n+1}; \hat{\theta}_n), i = 1, ..., n].$ 

#### Gaussian Process Model



Figure: Example of Gaussian Process model in 1-D

#### Challenge and Motivation

- Computational issue that hinders GP from broader application
  - Modelling and making inference involves manipulations of a  $n \times n$  correlation matrix  $R_n(\theta)$ , such as the calculation of  $R_n^{-1}(\theta)$  and  $|R_n(\theta)|$ . The computational order is  $\mathcal{O}(n^3)$ .
- The underestimation of GP predictor uncertainty
  - The resulting plug-in predictors tend to underestimate the uncertainty because variance is obtained by substituting the true parameters with their estimators.

#### LHD Example



Figure: Example of Latin Hypercube Design in 2-D

LHD-Based Block Bootstrap

#### Example of LHD-Based Block Bootstrap



Figure: d = 2, l = 24, b = 4, m = 6,  $|\mathcal{B}_n(\mathbf{i})| = 6$ , N = 36, n = 216

 LHD-Based Block Bootstrap

#### Example of LHD-Based Block Bootstrap



Figure: d = 2, l = 24, b = 4, m = 6,  $|B_n(i)| = 6$ , N = 36, n = 216

<ロ> < 部> < 書> < 言> こ 11/23

# Consistency of the Bootstrap Estimators: Converge in Probability

#### Theorem 1

Assume regularity conditions are satisfied. If  $m = o(n^{1/d})$  and  $m \to \infty$ , then

$$\hat{\phi}_{N}^{*}-\hat{\phi}_{n}
ightarrow 0$$
 prob  $-P_{N,\omega}^{*},$  prob  $-P.$ 

Note:  $\hat{T}_N^* \to 0 \quad \text{prob} - P_{N,\omega}^*, \text{prob} - P \text{ if for any } \epsilon > 0 \text{ and any } \delta > 0,$  $\lim_{n\to\infty} P\{P_{N,\omega}^*(|\hat{T}_N^* > \epsilon| > \delta)\} = 0.$ 

\* Yibo Zhao, Yasuo Amemiya and Ying Hung Efficient Gaussian Process Modelling Using Experimental Design-based Subagging *Statistica Sinaca* 

# **Objectives of Research**

Objective of this research is to construct a predictive distribution that is

- 1) easy to compute (due to the subsampling);
- 2) allow a better quantification of predictive uncertainty

#### Construction Methods

#### Definition 1 (Direct density prediction)

Given the realization  $\{X_n, y_n\}$ , let  $\{X_N^*, y_N^*\}$  be a bootstrap sample, a bootstrap predictive distribution is defined by

$$g^*(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n) = \int g(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_N^*, \boldsymbol{y}_N^*, \hat{\phi}_N^*) dP^*(\boldsymbol{X}_N^*, \boldsymbol{y}_N^* \mid \boldsymbol{X}_n, \boldsymbol{y}_n),$$

Based on the subsamples obtained from LHD-based bootstrap, a Monte Carlo estimate of the bootstrap predictive distribution can be obtained by

$$\tilde{g}^*(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n) = U^{-1} \sum_{u=1}^U g(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}^*_{N(u)}, \boldsymbol{Y}^*_{N(u)}, \hat{\phi}^*_{N(u)}),$$

When  $U \to \infty$ ,  $\tilde{g}^*(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n)$  converges to  $g^*(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n)$ .

## **Direct Density Prediction**



#### Construction Methods

#### Definition 2 (Normal approximation)

Utilizing LHD-based bootstrap approach, the Monte Carlo estimate of the predictive mean and variance are:

$$\begin{split} \tilde{\mu}^{*}(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_{n}, \boldsymbol{y}_{n}) &= \quad U^{-1} \sum_{u=1}^{U} \mu(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_{N(u)}, \boldsymbol{y}_{N(u)}, \hat{\phi}^{*}_{N(u)}) \\ \tilde{\sigma}^{2*}(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_{n}, \boldsymbol{y}_{n}) &= \quad U^{-1} \sum_{u=1}^{U} \sigma^{2}(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_{N(u)}, \boldsymbol{y}_{N(u)}, \hat{\phi}^{*}_{N(u)}). \end{split}$$

When  $U 
ightarrow \infty$ ,  $ilde{\mu}^*(m{x}_{n+1} \mid m{X}_n, m{y}_n)$  converges to

$$\mu^*(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n) = \int \mu(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_N^*, \boldsymbol{y}_N^*, \hat{\boldsymbol{\phi}}_N^*) dP^*(\boldsymbol{X}_N^*, \boldsymbol{y}_N^* \mid \boldsymbol{X}_n, \boldsymbol{y}_n)$$

and  $\tilde{\sigma}^{2*}(\pmb{x}_{n+1} \mid \pmb{X}_n, \pmb{y}_n)$  converges to

$$\sigma^{2*}(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_n, \boldsymbol{y}_n) = \int \sigma^2(\boldsymbol{x}_{n+1} \mid \boldsymbol{X}_N^*, \boldsymbol{y}_N^*, \hat{\boldsymbol{\phi}}_N^*) dP^*(\boldsymbol{X}_N^*, \boldsymbol{y}_N^* \mid \boldsymbol{X}_n, \boldsymbol{y}_n).$$

# Normal Approximation Prediction



#### Theoretical Comparison

#### Theorem 2

Let  $\sum_{i}$  be the summation of all  $m^{d}$  blocks and  $\sum_{\pi_{1},...,\pi_{d}}$  be the summation of independent permutation over  $\{0, 1, ..., m-1\}$ . Under assumption (A.3), we have

• Direct density and normal approximation both have **unbiased** predictive mean. i.e.,

$$E(\mu(\mathbf{x}_{n+1} \mid \mathbf{X}_n, \mathbf{y}_n, \hat{\phi}_n) - \mu_1^*) = E(\mu(\mathbf{x}_{n+1} \mid \mathbf{X}_n, \mathbf{y}_n, \hat{\phi}_n) - \mu_2^*)$$
  
=  $E(\frac{m^{d-1} - 1}{m^{d-1}} \sum_{i} \gamma_i(\hat{\theta}_n)^T R_{i,i}^{-1}(\hat{\theta}_n)(\mathbf{y}_i - \mathbf{X}_i \hat{\beta}_n) + O(N^{-1/2}))$   
= 0

#### Simulation Setting

• 
$$Y(\mathbf{x}) = \mu(\mathbf{x}) + Z(\mathbf{x})$$

• 
$$\mu(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$$

• 
$$\psi(\mathbf{x}_1 - \mathbf{x}_2) = \exp(-\sum_{i=1}^2 |x_{1i} - x_{2i}|/\theta_i)$$

• 
$$\beta_1 = \beta_2 = \theta_1 = \theta_2 = \sigma^2 = 1$$

- Generate n=400, 2500 realizations on regular grid  $[0,1]^2$
- For each choice of sample size, a total of 100 data sets are simulated
- 20 LHD-based block bootstrap samples are collected

#### Simulation studies

Table: Comparisons of prediction in three untried settings with 100 replications (standard deviation in parenthesis).

Method	Summary statistics	$x_{n+1}$	$x_{n+2}$	<b>x</b> <sub>n+3</sub>
		<i>n</i> = 400		
Plug-in	Mean	0.59 (0.9965)	1.05 (1.0447)	1.10 (1.0558)
	Variance	0.01 (0.0030)	0.03 (0.0061)	0.02 (0.0040)
Density $m = 4$	Mean	0.63 (0.8666)	1.05 (0.9383)	1.09 (0.9533)
	Variance	0.22 (0.0559)	0.15 (0.0298)	0.14 (0.0291)
Density $m = 6$	Mean	0.62 (0.8416)	1.08 (0.9604)	1.12 (0.9683)
	Variance	0.25 (0.0741)	0.15 (0.0388)	0.13 (0.0380)
Normal $m = 4$	Mean	0.63 (0.8666)	1.05 (0.9383)	1.09 (0.9533)
	Variance	0.17 (0.0287)	0.11 (0.0173)	0.09 (0.0150)
Normal $m = 6$	Mean	0.62 (0.8416)	1.08 (0.9604)	1.12 (0.9683)
	Variance	0.14 (0.0249)	0.11 (0.0165)	0.10 (0.0148)
		<i>n</i> = 2500		
Plug-in	Mean	0.67 (0.9785)	1.10 (1.1682)	1.14 (1.1731)
	Variance	0.02 (0.0017)	0.02 (0.0016)	0.01 (0.0013)
Density $m = 4$	Mean	0.66 (0.9183)	1.12 (1.0952)	1.14 (1.0899)
	Variance	0.21 (0.0517)	0.11 (0.0272)	0.10 (0.0254)
Density $m = 6$	Mean	0.62 (0.8931)	1.10 (1.0751)	1.13 (1.0710)
	Variance	0.19 (0.0401)	0.11 (0.0204)	0.11 (0.0216)
Normal $m = 4$	Mean	0.66 (0.9183)	1.12 (1.0952)	1.14 (1.0899)
	Variance	0.15 (0.0129)	0.07 (0.0057)	0.06 (0.0050)
Normal $m = 6$	Mean	0.62 (0.8931)	1.10 (1.0751)	1.13 (1.0710)
	Variance	0.13 (0.0115)	0.08 (0.0065)	0.08 (0.0061)

・ロ・・ (日・・ヨ・・ヨ・・ ヨー・ つう

# Ongoing Work

- Compare the predictive variance of direct density prediction and normal approximation with the case when full data is used.
- Quantify the gain in predictive variance of LHD-based block bootstrap with SRS block bootstrap.

## Summary

- LHD-based block bootstrap borrows the strength of space-filling designs to provide an efficient subsampling plan and reduce computational complexity.
- Two methods are proposed to construct **bootstrap predictive distributions**.
- We show the **unbiasedness** of the predictive mean under both direct density prediction and normal approximation prediction.

# Thank you!