

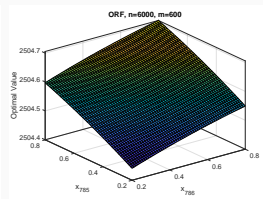
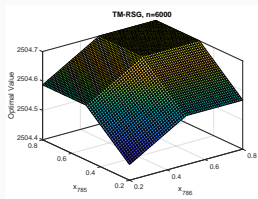
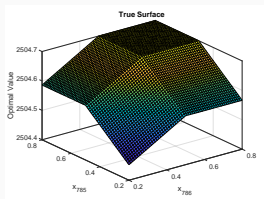
Sample and Computationally Efficient Simulation Metamodeling in High Dimensions

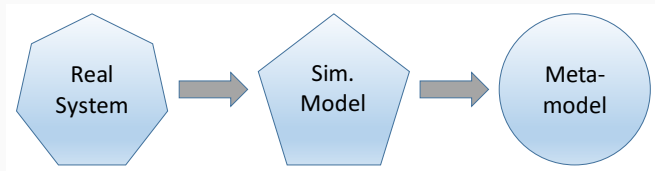
Xiaowei Zhang (HKU Business School)

December 5, 2020

Joint work with Liang Ding (Texas A&M University)

Surface Fitting in 816 Dimensions





- Simulation models are often computationally expensive
- Metamodel: **statistical** model for simulation input-output relationship
 - A.K.A. surrogate model
 - Run simulation at a small number of design points
 - Predict responses with the fitted statistical model

Comparison with Regression

	Metamodeling	Regression
Statistical Model	Linear/Nonparametric	Linear/Nonparametric
Typical Design	Fixed	Random
Typical Noise	Heteroscedastic	Homoscedastic
Replications	Yes	No

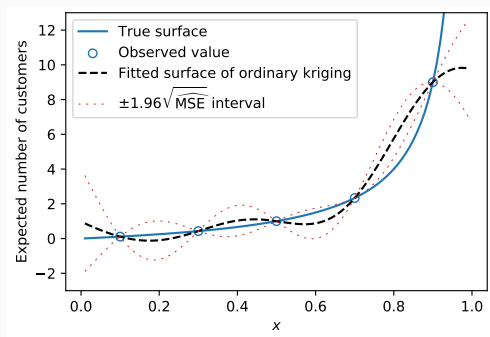
Stochastic Kriging

- Gaussian process regression
- Response surface is a sample path of a GP with kernel $k(x, x')$

$$Y(\cdot) \sim \text{GP}(0, k(\cdot, \cdot))$$

- Take samples at design points $\{x_1, \dots, x_n\}$
- SK predictor

$$\hat{Y}_n(x) = \mathbf{k}^T(x)(\mathbf{K} + \Sigma)^{-1}\bar{\mathbf{Y}}$$



The prediction accuracy of SK depends on

- (i) choice of kernel
- (ii) choice of experimental design

- Gaussian kernel

$$k_{\text{Gauss}}(\mathbf{x}, \mathbf{x}') := \exp(-\|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\|^2)$$

- Matérn kernel

$$k_{\text{Matérn}(\nu)}(\mathbf{x}, \mathbf{x}') := \frac{1}{2^{\nu-1}\Gamma(\nu)} \left(\sqrt{2\nu} \|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\| \right)^\nu K_\nu(\sqrt{2\nu} \|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\|)$$

- Gaussian kernel

$$k_{\text{Gauss}}(\mathbf{x}, \mathbf{x}') := \exp(-\|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\|^2)$$

- Matérn kernel

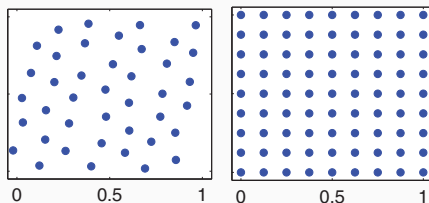
$$k_{\text{Matérn}(\nu)}(\mathbf{x}, \mathbf{x}') := \frac{1}{2^{\nu-1}\Gamma(\nu)} \left(\sqrt{2\nu} \|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\| \right)^\nu K_\nu(\sqrt{2\nu} \|\boldsymbol{\theta}^\top(\mathbf{x} - \mathbf{x}')\|)$$

- Generalized integrated Brownian field (Salemi et al. 2019, *OR*)
 - Motivation: model smoothness separately for each dimension

$$k(\mathbf{x}, \mathbf{x}') = \prod_{j=1}^d k_j(x_j, x'_j)$$

- On dimension j is a Brownian motion integrated ℓ_j times

Typical Experimental Designs



- Latin hypercube design (LHD)
 - Pros: ease of use; 1-d projections onto are evenly dispersed
 - Con: superiority only holds for large sample sizes when d is large
- Lattice design: Cartesian product of one-dimensional designs
 - Pro: $\mathbf{K}^{-1} = \bigotimes_{j=1}^d \mathbf{K}_j^{-1}$ for tensor-product kernels
 - Con: excessive sample size when design space when d is large

- **Sample** complexity: n grows exponentially in d
- **Computational** complexity: $\mathcal{O}(n^3)$

- Sample paths of a GP with Matérn(α) kernel form a *Sobolev space*
 - α represents the smoothness
- Minimax-optimal rate for estimating Y via noisy samples: $n^{-\alpha/(2\alpha+d)}$
- So, sample complexity for achieving an δ -error is $\delta^{-(2+d/\alpha)}$
- Practical situation could be even worse due to model misspecification
 - SK is specified with Matérn(ν) kernel
 - Convergence rate: $\mathcal{O}(n^{-\min(\alpha,\nu)/(2\nu+d)})$

Combat the Curse with Smoothness?

- Given a large d , the convergence rate is fast with a large α
- Why don't we set/assume $\alpha = \infty$?
 - Infinitely differentiable response surface is rare
 - E.g., the max function often appears in queueing, inventory, FE models

- Computing $(\mathbf{K} + \Sigma)^{-1}$ requires $\mathcal{O}(n^3)$
- Subsampling to construct a low-rank approximation: $\mathcal{O}(\ell^2 n)$
- Lu et al. (2020, *OR*), but huge literature in machine learning and statistics

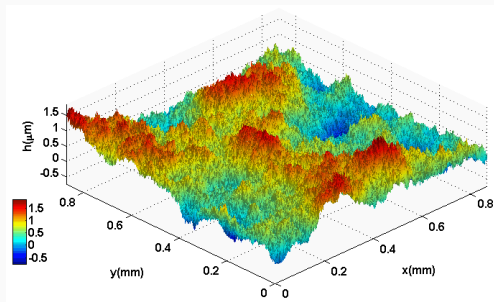
- “New” kernel: Tensor Markov (TM) kernels
- New exponential design: Random sparse grid (RSG) designs
- Convergence rate: “weakly” dependent of d
- Fast, **exact** computation

- Tensor-product form: $k(\mathbf{x}, \mathbf{x}') = \prod_{j=1}^d k_j(x_j, x'_j)$
- Each k_j corresponds to a Gauss-Markov process
 - Brownian motion: $k(x, x') = x \wedge x'$
 - Stationary OU process: $k(x, x') = \exp(-\theta|x - x'|)$

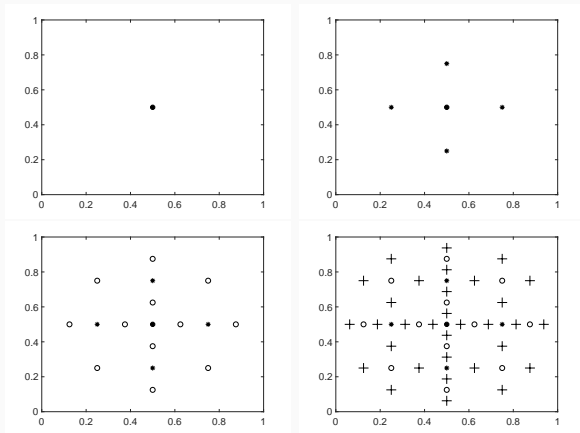
- Facilitate computation
- One-dimensional case: K^{-1} is **tri-diagonal**
 - Assume $\mathcal{X} = [0, 1]$ and $x_i = \frac{i}{n+1}$, $i = 1, \dots, n$
 - Brownian motion: $k(x, y) = \min(x, y)$

$$K^{-1} = (n+1) \begin{pmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & \ddots & \ddots & \ddots & & \\ & & -1 & 2 & -1 & \\ & & & -1 & 1 & \end{pmatrix}$$

- Time-changed Brownian field



Classical Sparse Grids



- Controlled by a **level** parameter τ

Dimension d	Full Grid	Sparse Grid of Level 4
1	15	15
2	225	49
5	759,375	351
10	5.77×10^{11}	2,001
20	3.33×10^{23}	13,201
50	6.38×10^{58}	182,001

- Pro: Fast computation of K^{-1} for tensor-product kernels (Plumlee, 2014, *JASA*)

Pro and Con of Classical Sparse Grids

- Pro: Fast computation of K^{-1} for tensor-product kernels (Plumlee, 2014, JASA)

Level τ	$d = 2$	$d = 5$	$d = 10$	$d = 20$	$d = 50$
2	5	11	21	41	101
3	17	71	241	881	5,201
4	49	351	2,001	13,201	182,001
5	129	1,471	13,441	154,881	4,867,201

- Con: Inflexible to use, fast computation only available for **complete** SGs
 - n must be coincide with the sample size of some level τ

- Find τ such that n falls between $\mathcal{X}_\tau^{\text{SG}}$ and $\mathcal{X}_{\tau+1}^{\text{SG}}$
- Random sampling on $\mathcal{X}_{\tau+1}^{\text{SG}} \setminus \mathcal{X}_\tau^{\text{SG}}$
- $\mathcal{X}_n^{\text{RSG}} := \mathcal{X}_\tau^{\text{SG}} \cup \mathcal{A}$
- Fast computation of \mathbf{K}^{-1} only for TM kernels

- Model well-specified: true surface is a GP with **tensor Markov** kernel

$$\max_{\mathbf{x} \in [0,1]^d} \mathbb{E}[(\hat{Y}_n(\mathbf{x}) - Y(\mathbf{x}))^2] = \mathcal{O}\left(n^{-1}(\log n)^{2(d-1)}\right)$$

- Model well-specified: true surface is a GP with **tensor Markov** kernel

$$\max_{\mathbf{x} \in [0,1]^d} \mathbb{E}[(\hat{Y}_n(\mathbf{x}) - Y(\mathbf{x}))^2] = \mathcal{O}\left(n^{-1}(\log n)^{2(d-1)}\right)$$

- If true surface is a GP with Matérn(α) and SK uses the same kernel, then

$$\max_{\mathbf{x} \in [0,1]^d} \mathbb{E}[(\hat{Y}_n(\mathbf{x}) - Y(\mathbf{x}))^2] = \mathcal{O}\left(n^{-2\alpha/(2\alpha+d)}\right).$$

- Smooth surface, rough kernel
- True surface is a GP with a **tensor product** kernel that is **smoother**

$$\max_{\mathbf{x} \in [0,1]^d} \mathbb{E}[(\hat{Y}_n^{\text{mis}}(\mathbf{x}) - Y^*(\mathbf{x}))^2] = \mathcal{O}\left(n^{-2}(\log n)^{3(d-1)}\right).$$

- Model well-specified:

$$\mathcal{O} \left(n^{-1} (\log n)^{2(d-1)} + (\log n)^d \max_{1 \leq i \leq n} m_i^{-1/2} \sigma(\mathbf{x}_i) \right).$$

- Model mis-specified:

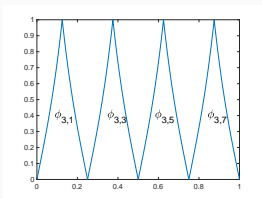
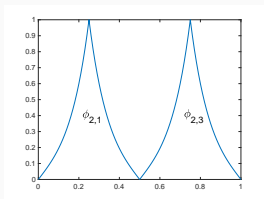
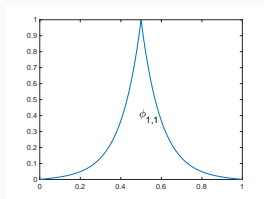
$$\mathcal{O} \left(n^{-2} (\log n)^{3(d-1)} + (\log n)^d \max_{1 \leq i \leq n} m_i^{-1/2} \sigma(\mathbf{x}_i) \right).$$

Key for Proof: Orthogonal Expansions

- Expansion for TM kernels

$$k(x, x') = \sum_{\tau=1}^{\infty} \sum_{(l,i): c_{l,i} \in \mathcal{X}_{\tau}^{\text{SG}}} \frac{\phi_{l,i}(x)\phi_{l,i}(x')}{\|\phi_{l,i}\|_{\mathcal{H}_k}^2}$$

- $\phi_{l,i} \in [0, 1]$
- $\|\phi_{l,i}\|_{\mathcal{H}_k}^2 \asymp 2^{|l|}$



- Expansion for GP with TM kernels

$$Y(\mathbf{x}) = \sum_{\tau=1}^{\infty} \sum_{(l,i): \mathbf{c}_{l,i} \in \mathcal{X}_{\tau}^{\text{SG}}} \frac{\phi_{l,i}(\mathbf{x})}{\|\phi_{l,i}\|_{\mathcal{H}_k}} Z_{l,i}$$

- Observing $\{Y(\mathbf{c}_{l,i}) : \mathbf{c}_{l,i} \in \mathcal{X}_{\tau}^{\text{SG}}\}$ equals observing $\{Z_{l,i} : \mathbf{c}_{l,i} \in \mathcal{X}_{\tau}^{\text{SG}}\}$

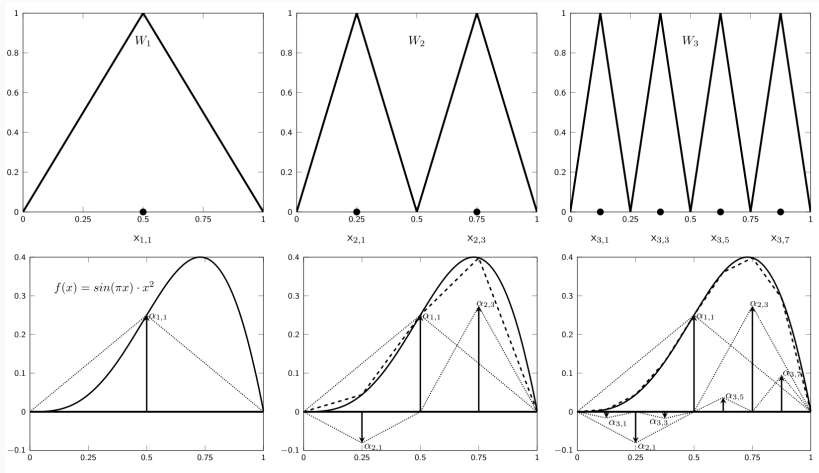


Figure 1: Brumm and Scheidegger (2017, *Econometrica*)

$$\begin{aligned}
\mathbb{E}[(\hat{Y}_n(\mathbf{x}) - Y(\mathbf{x}))^2] &= \sum_{\tau=\ell+1}^{\infty} \sum_{(\mathbf{l}, \mathbf{i}): \mathbf{c}_{\mathbf{l}, \mathbf{i}} \in \mathcal{X}_{\tau}^{\text{SG}}} \frac{\phi_{\mathbf{l}, \mathbf{i}}^2(\mathbf{x})}{\|\phi_{\mathbf{l}, \mathbf{i}}\|_{\mathfrak{H}_{\mathcal{C}_k}}^2} \\
&\leq \sum_{\tau=\ell+1}^{\infty} \sum_{(\mathbf{l}, \mathbf{i}): \mathbf{c}_{\mathbf{l}, \mathbf{i}} \in \mathcal{X}_{\tau}^{\text{SG}}} \frac{1}{\|\phi_{\mathbf{l}, \mathbf{i}}\|_{\mathfrak{H}_{\mathcal{C}_k}}^2} \\
&\asymp \sum_{|\mathbf{l}| > \ell + d - 1} 2^{-|\mathbf{l}|} \\
&= \mathcal{O}(2^{-\ell} \ell^{d-1}) \\
&= \mathcal{O}\left(n^{-1} (\log n)^{2(d-1)}\right)
\end{aligned}$$

- K^{-1} can be expressed as

$$K^{-1} = \begin{pmatrix} & |\mathcal{X}_\tau^{\text{SG}}| \text{ dim.} & \tilde{n} \text{ dim.} \\ \square & & \square \\ \square & & \mathbf{D} \end{pmatrix}$$

- Each block can be computed efficiently
- K^{-1} is **sparse**: Proportion of nonzero entries: $\mathcal{O}(n^{-1}(\log n)^{2d})$

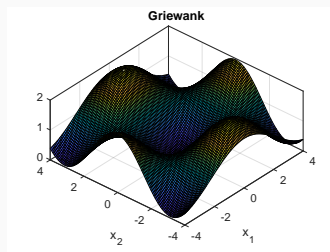
- \mathbf{K}^{-1} can be expressed as

$$\mathbf{K}^{-1} = \begin{pmatrix} & |\mathcal{X}_\tau^{\text{SG}}| \text{ dim.} & \tilde{n} \text{ dim.} \\ \square & & \square \\ \square & & \mathbf{D} \end{pmatrix}$$

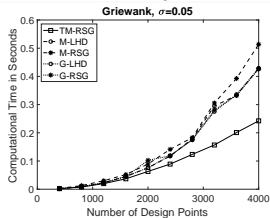
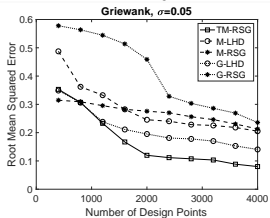
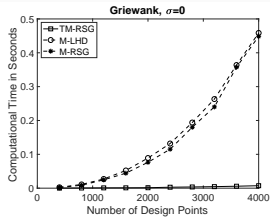
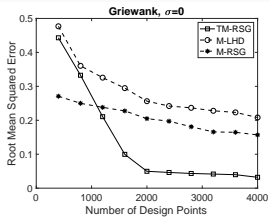
- Each block can be computed efficiently
- \mathbf{K}^{-1} is **sparse**: Proportion of nonzero entries: $\mathcal{O}(n^{-1}(\log n)^{2d})$
- Computing stochastic kriging: Woodbury matrix identity

$$(\mathbf{K} + \Sigma)^{-1} = \Sigma^{-1} - \Sigma^{-1}(\mathbf{K}^{-1} + \Sigma^{-1})^{-1}\Sigma^{-1},$$

Griewank Function in 10 Dimensions



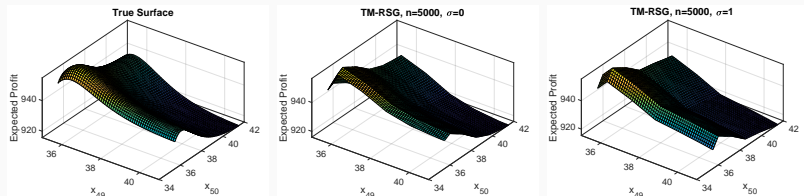
$$y_{\text{Griewank}}(\mathbf{x}) = \sum_{j=1}^d \frac{x_j^2}{4000} - \prod_{j=1}^d \cos\left(\frac{x_j}{\sqrt{j}}\right) + 1, \quad \mathbf{x} \in [-4, 4]^d, \quad d = 10$$

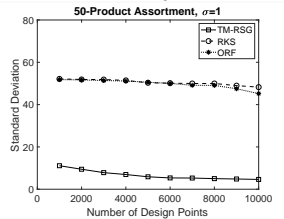
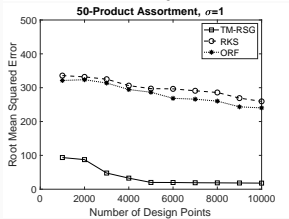
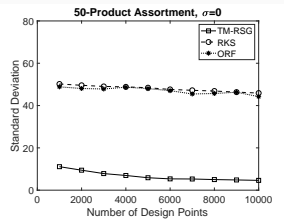
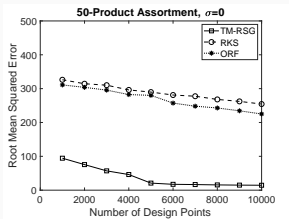


A Product Assortment Problem

- Aydin and Porteus (2008, *OR*)
- Design variable ($d = 50$): price vector
- Response surface: expected profit

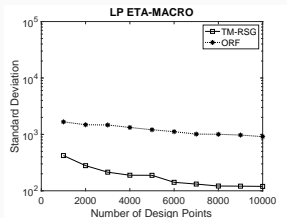
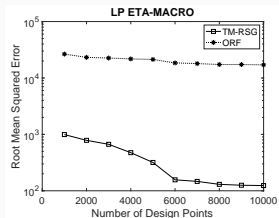
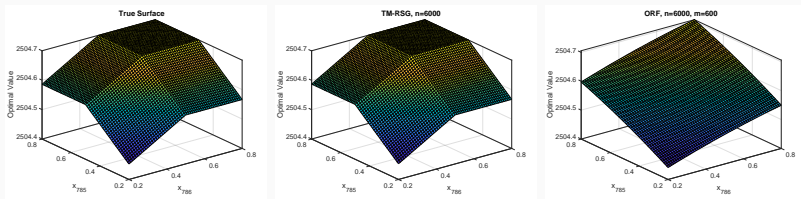
$$y(\mathbf{x}) = \frac{1}{2(b-a)} \sum_{j=1}^d \left[(b-a) \left(\frac{x_j - c_j}{x_j} \right) + a \right]^2 Q_j^2(\mathbf{x})$$





A Large Linear Program

- Decision variable ($d = 816$): coefficient vector of the objective function
- Response surface: optimal value



- Curse of dimensionality
- Tensor Markov kernels
- Random sparse grids
- **Sample** efficiency: convergence rate suffers little from curse of dimensionality
- **Computational** efficiency: exact computation from sparse structure of K^{-1}