

# Kriging and Alternatives in Computer Experiments

C. F. Jeff Wu

ISyE, Georgia Institute of Technology

- Use kriging to build meta-models in computer experiments, a brief review
- Numerical problems with kriging
- Alternatives to kriging:
  - Regularized kriging, Hybrid kriging
  - Overcomplete basis surrogate model (OBSM)

## Why computer experiments?

- ✓ No need for expensive lab equipments and materials, less costly than physical experiments.



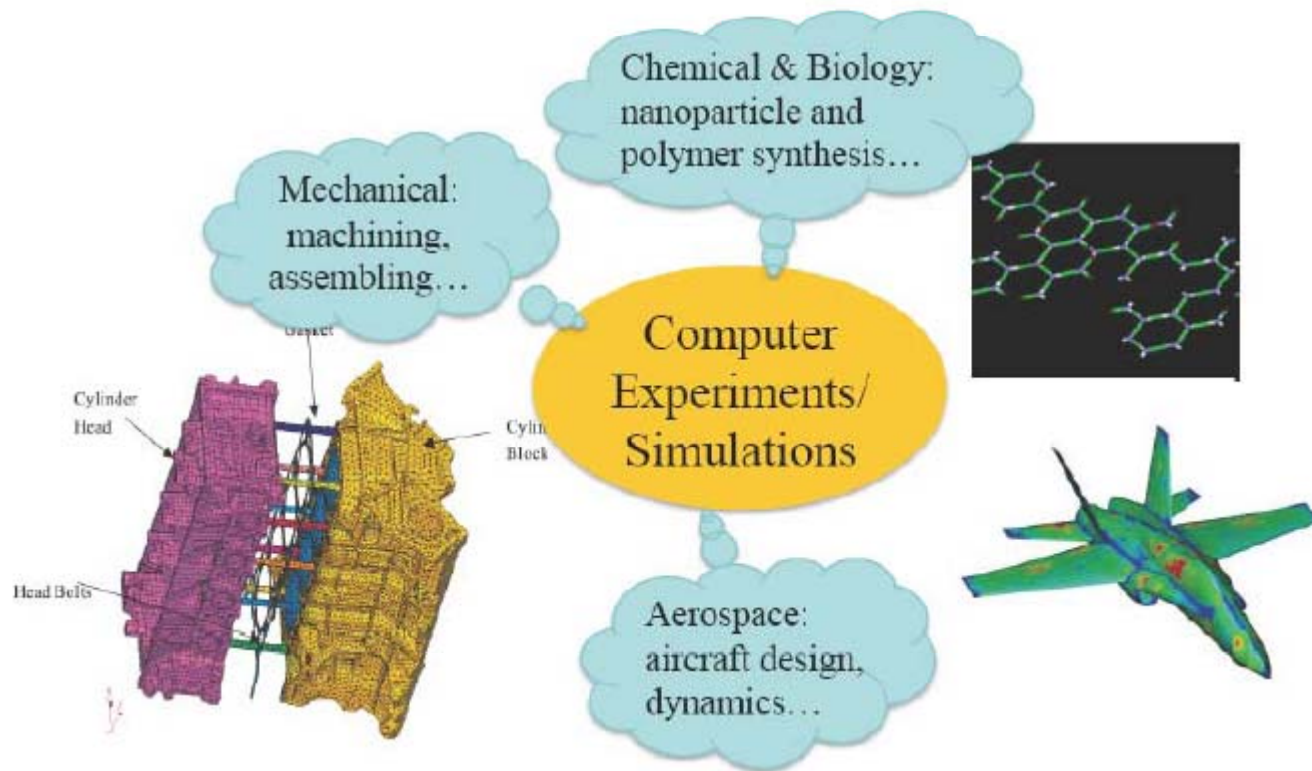
- ✓ Not affected by human and environmental factors.



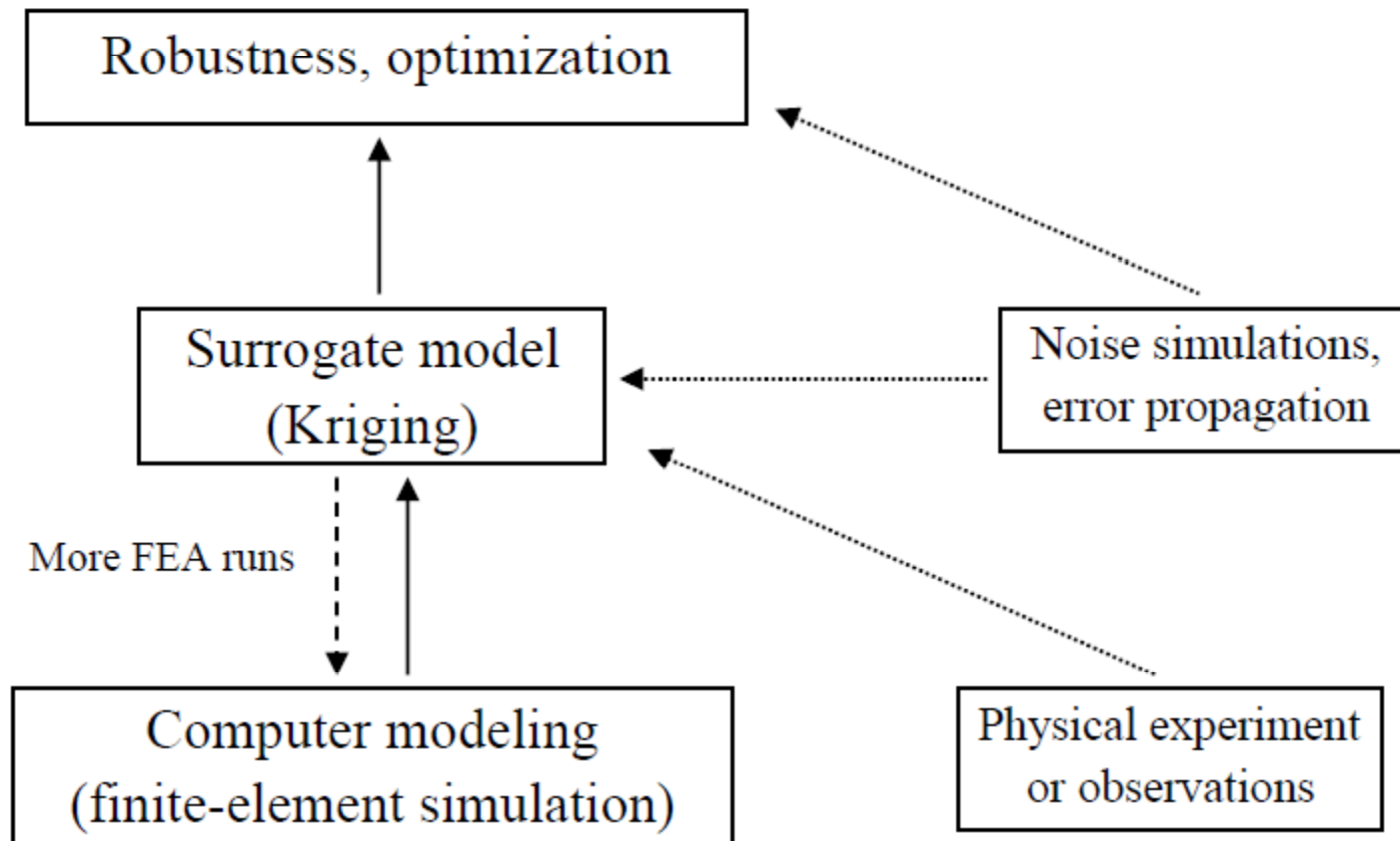
- ✓ Study dangerous or infeasible physical experiments, such as ammunition detonation.



## Some examples



## Statistical Meta-Modeling of Computer Experiments



# Kriging Models

- Ordinary Kriging

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

$$Z(\mathbf{x}) \sim N_n(\mathbf{0}, \sigma^2 \varphi(\mathbf{h})) \equiv GP(\mathbf{0}, \sigma^2 \varphi(\mathbf{h}))$$

- Correlation function

- $\varphi(\mathbf{0}) = 1$
- $\varphi(\mathbf{h}) = \varphi(-\mathbf{h})$ , (symmetric function)
- $\varphi$  is a positive semi-definite function

# Correlation function

- Matérn

$$\varphi(h) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(2\sqrt{\nu}\theta|h|\right)^\nu K_\nu\left(2\sqrt{\nu}\theta|h|\right)$$

where  $K_\nu$  is the modified Bessel function of order  $\nu$

$$\nu \rightarrow \infty, \quad \varphi(h) \rightarrow \exp(-\theta h^2)$$

- *Power exponential correlation*

$$\varphi(h) = \exp\left(-\theta|h|^q\right), \quad 0 < q \leq 2, \quad 0 < \theta$$

- $q=2$  *Gaussian correlation function* (infinitely differentiable)
- $q=1$  Ornstein-Uhlenbeck process ( $\nu=1$  in Matérn)

- Linear, Cubic correlation

# Kriging predictor

- Best Linear Unbiased Predictor (BLUP)

$$\hat{y}(\mathbf{x}) = \hat{\mu} + r(\mathbf{x})' \mathbf{R}^{-1}(\mathbf{y} - \hat{\mu} \mathbf{1}),$$


$$r(\mathbf{x})' = (\varphi(\mathbf{x} - \mathbf{x}_1), \dots, \varphi(\mathbf{x} - \mathbf{x}_n)),$$

$$\hat{\mu} = \mathbf{1}' \mathbf{R}^{-1} \mathbf{y} / \mathbf{1}' \mathbf{R}^{-1} \mathbf{1},$$


$$\hat{y}(\mathbf{x}_i) = y_i \quad \text{an interpolating property}^*$$

\*required for deterministic simulations

# Recent work in kriging

- Calibration of computer model, Kriging with *calibration* parameters (Kennedy-O'Hagan, 2001), with *tuning* parameters (Santner et al., 2009)
- Computer simulations with **different** levels of accuracy (Kennedy-O'Hagan, 2000; Qian et al., 2006; Qian-Wu, 2008)   
construction of **nested** space-filling (e.g., Latin hypercube) designs (Qian-Ai-Wu, 2009, various papers by Qian and others, 2009-date)

# Recent work in kriging (cont.)

- Kriging for multiple outputs and functional response (Conti et al., 2009; Conti and O'Hagan, 2010)
- Treed Gaussian Process model (Gramacy and Lee, 2008).
- Kriging (i.e., GP model) with quantitative and **qualitative** factors (Qian-Wu-Wu, 2008, Han et al., 2009)   
construction of **sliced** space-filling (e.g., Latin hypercube) designs (Qian-Wu, 2009, Qian, 2010)

# Maximum Likelihood Estimation

- Profile log-likelihood approach

$$Q(\boldsymbol{\theta}) = n \log(\sigma^2(\boldsymbol{\theta})) + \log|\mathbf{R}(\boldsymbol{\theta})|$$

where  $\sigma^2(\boldsymbol{\theta}) = \{\mathbf{y} - \hat{\mu}(\boldsymbol{\theta})\mathbf{1}\}'\mathbf{R}^{-1}(\boldsymbol{\theta})\{\mathbf{y} - \hat{\mu}(\boldsymbol{\theta})\mathbf{1}\}/n$

$$\hat{\mu}(\boldsymbol{\theta}) = \mathbf{1}'\mathbf{R}^{-1}(\boldsymbol{\theta})\mathbf{y} / \mathbf{1}'\mathbf{R}^{-1}(\boldsymbol{\theta})\mathbf{1}$$

# Numerical Instability in $R^{-1}(\theta)$

- $R(\theta)$  is an  $n \times n$  matrix,  $n$ =sample size
- Its condition number (max e.v./min e.v.)  $\uparrow$  as

I. Sample size  $n \uparrow$

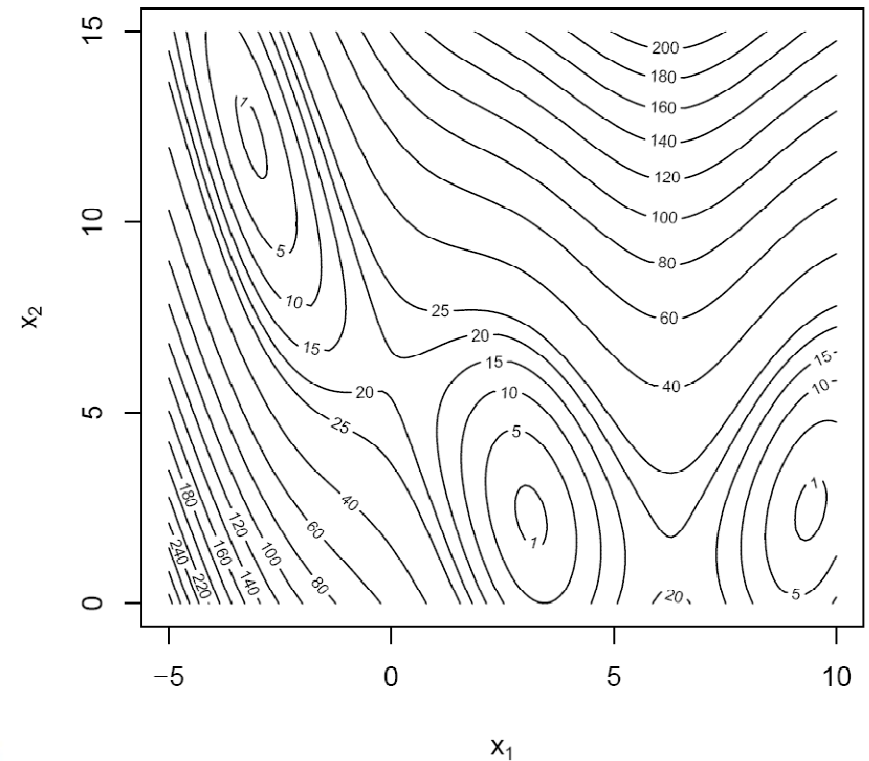
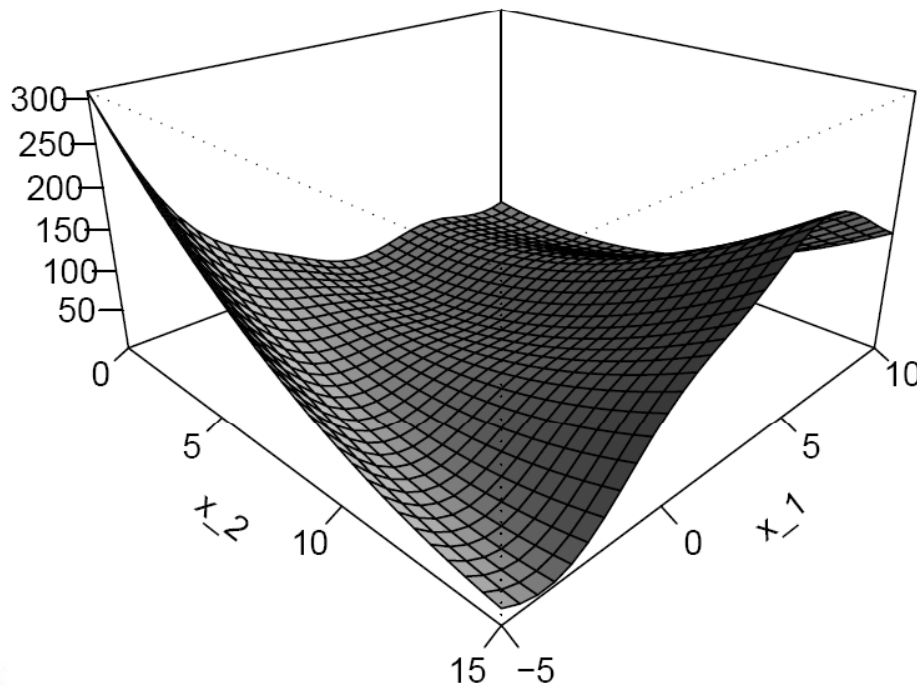
II. Dimension of input vectors  $\uparrow$

(Peng-Wu, 2010)

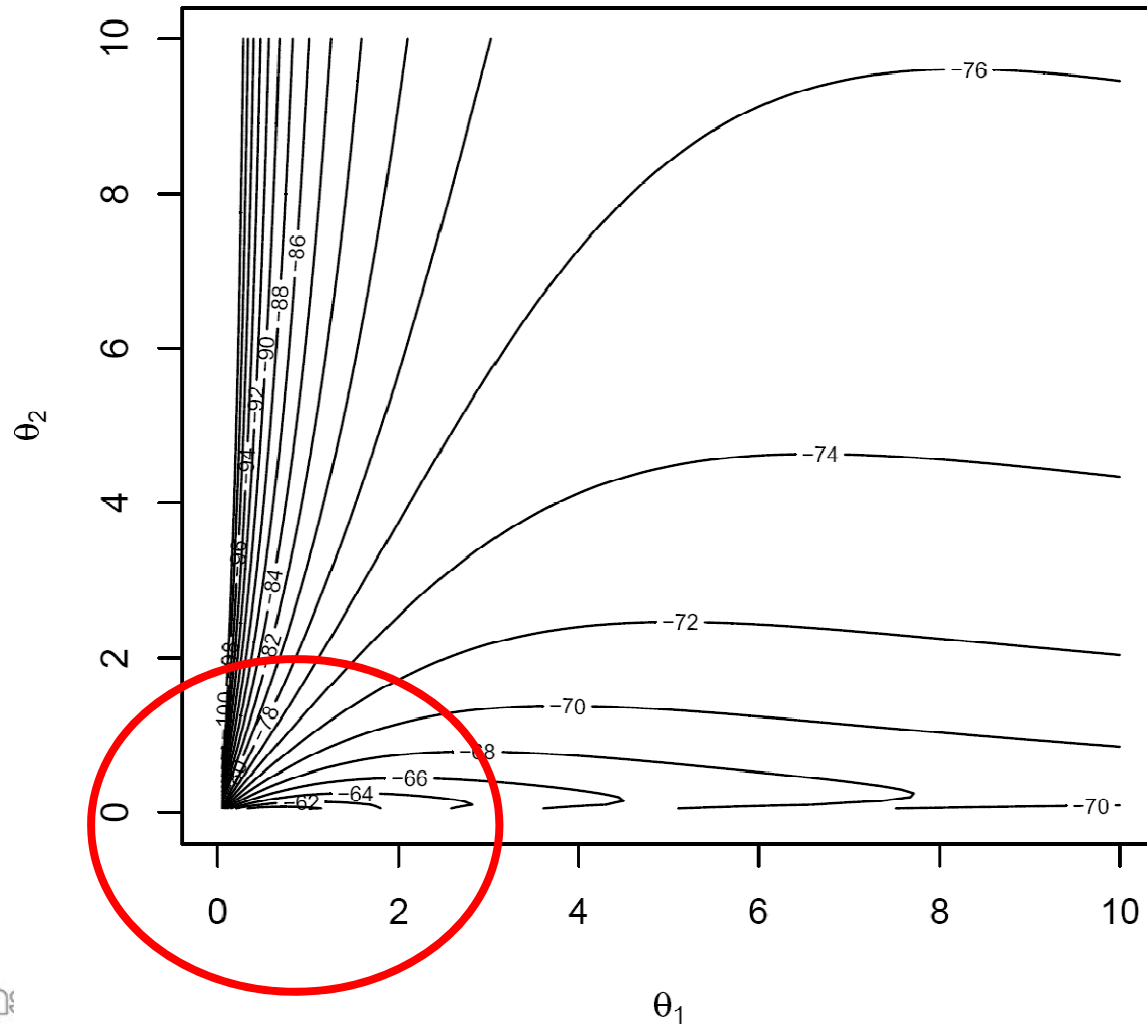
# Branin function

(Andre, Siarry and Dognon, 2001)

$$f(x_1, x_2) = \left( x_2 - \frac{5.1}{4\pi^2} x_1^2 + \frac{5}{\pi} x_1 - 6 \right)^2 + 10 \left( 1 - \frac{1}{8\pi} \right) \cos(x_1) + 10$$



# Log-likelihood function (Regular grid: $n = 7^2$ )



# Regularized Kriging

- Introducing a regularizing constant  $\lambda$  into the predictor

$$\hat{y}_\lambda(\mathbf{x}) = \hat{\mu}_\lambda + r(\mathbf{x})'(\mathbf{R} + \lambda\mathbf{I})^{-1}(\mathbf{y} - \hat{\mu}_\lambda\mathbf{1})$$

$$\text{where } \hat{\mu}_\lambda = \mathbf{1}'(\mathbf{R} + \lambda\mathbf{I})^{-1}\mathbf{y} / \mathbf{1}'(\mathbf{R} + \lambda\mathbf{I})^{-1}\mathbf{1}$$

Peng and Wu (2010, submitted)

- Similar modification in estimation: maximizing a **regularized likelihood**

# Kriging with nugget effects

- Model from spatial statistics

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

- BLUP

$$\hat{y}_\delta(\mathbf{x}) = \hat{\mu}_\delta + r(\mathbf{x})' \left( \mathbf{R} + \frac{\delta^2}{\sigma^2} \mathbf{I} \right)^{-1} (\mathbf{y} - \hat{\mu}_\delta \mathbf{1})$$

$$\hat{\mu}_\delta = \frac{\mathbf{1}' \left( \mathbf{R} + \frac{\delta^2}{\sigma^2} \mathbf{I} \right)^{-1} \mathbf{y}}{\mathbf{1}' \left( \mathbf{R} + \frac{\delta^2}{\sigma^2} \mathbf{I} \right)^{-1} \mathbf{1}}$$

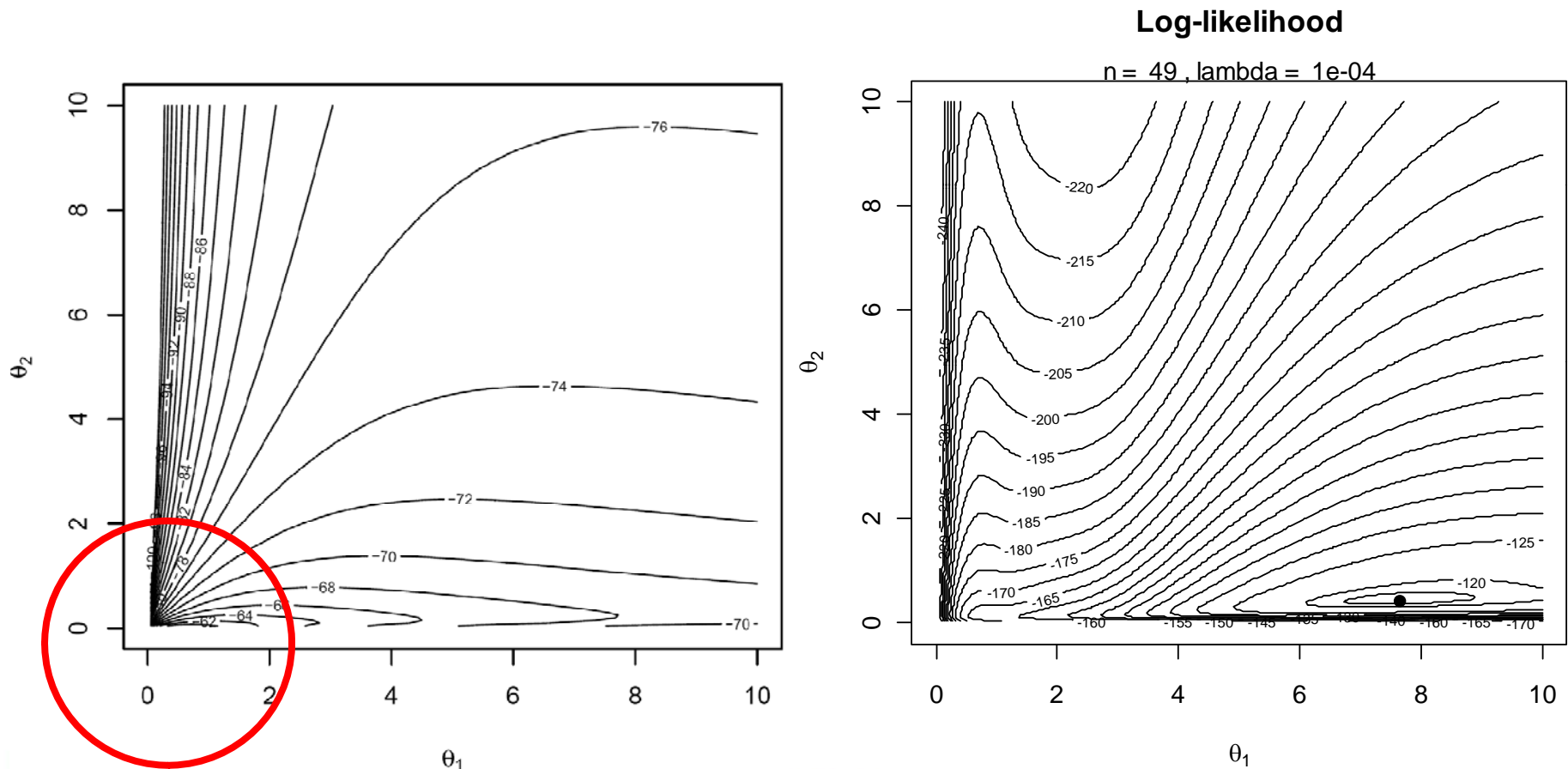
# Algorithm (Ridge Trace)

- Root Mean Squared Prediction Error (RMSPE)

$$\text{RMSPE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (Y(\mathbf{x}_i) - \hat{Y}(\mathbf{x}_i))^2}$$

- (1) Set  $\lambda^*$  as the lower bound and choose a grid point set for  $\lambda$ , say,  $(\lambda_1, \dots, \lambda_k)$ , and let  $i = 1$ .
- (2) Use  $\lambda_i$  in regularized kriging to estimate  $\theta_{\lambda}$ .
- (3) Compute the RMSPE. Let  $i = i + 1$ .
- (4) Repeat steps 2 and 3 until all  $k$  grid points are exhausted.
- (5) The final estimator  $\hat{\theta}_{\hat{\lambda}}$  is the one with the lowest RMSPE with  $\hat{\lambda}$ .

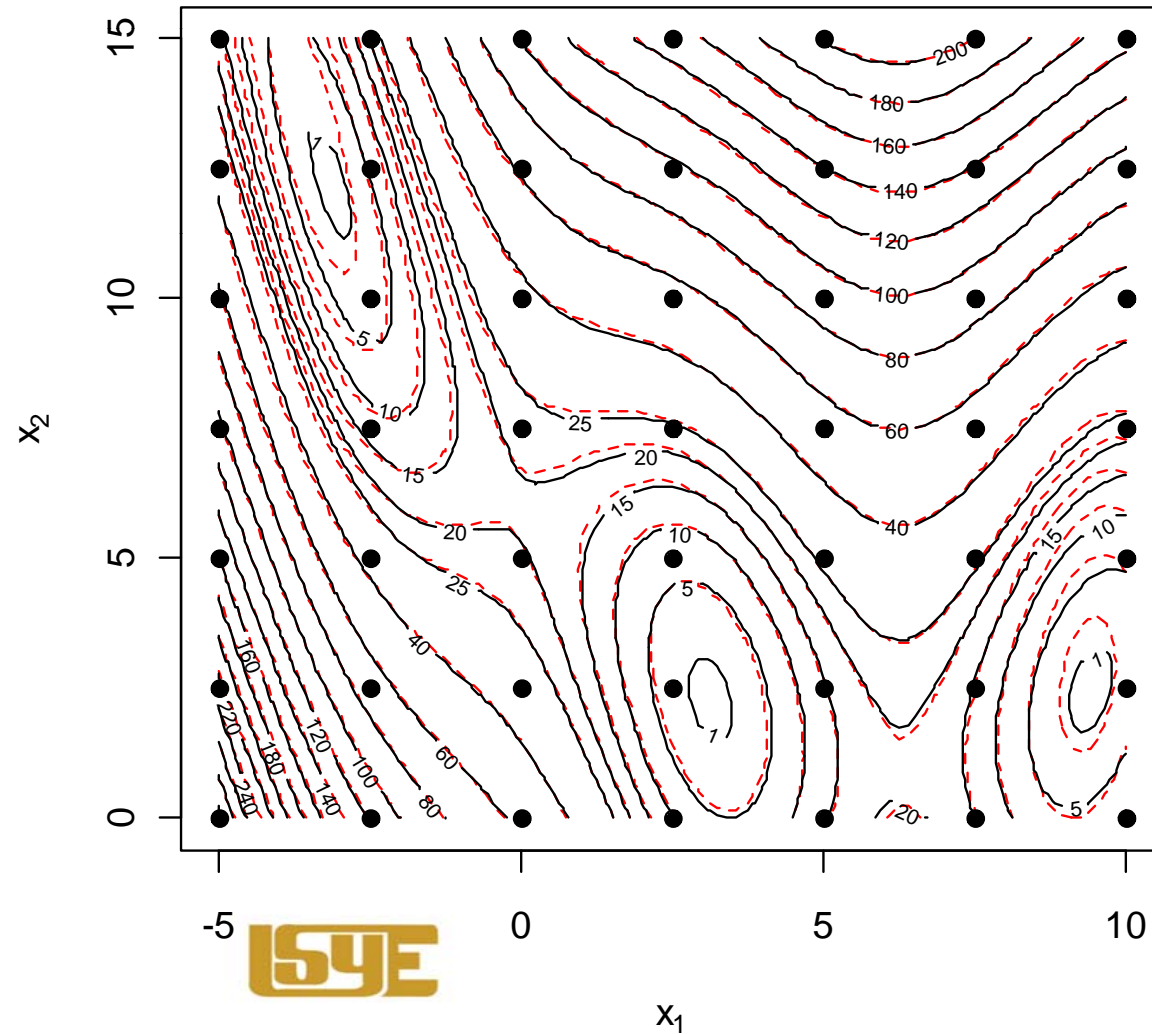
# Log-likelihood function, Branin function, (Regular grid: $n = 7^2$ )



$$\lambda^* = \frac{n^{1/2}}{\Delta^{1/2} - 1} = 3.3 \times 10^{-7}$$

## Regularized Kriging

$\lambda$	RMSPE
$10^{-1}$	6.1842
$10^{-2}$	3.2085
$10^{-3}$	1.5605
$10^{-4}$	1.0706
$10^{-5}$	1.7069
$10^{-6}$	2.5104
$10^{-7}$	3.2348



# Overcomplete Basis Surrogate Model

- Use an overcomplete dictionary of basis functions
- Use linear combinations of basis functions to approximate unknown functions
- Use Matching Pursuit for fast (i.e. greedy) computations
- Choice of basis functions to “mimic” the shape of the surface

Chen, Wang, and Wu (2010, *IIE Tran.* Q&R)

# Surrogate Representation

- Surrogate model: use a linear combination of pre-specified basis functions, i.e.,

$$f(\mathbf{x}) = \sum_j c_j \phi_j(\mathbf{x}), \mathbf{x} \in \mathcal{X}$$

- unitary norm  $\|\phi_j\| = 1$
- basis dictionary,  $\{\phi_j, j = 1, \dots, M\}$
- no unknown parameter in  $\phi_j$ , only unknown are the *linear*  $c_j$
- Overcomplete :  $M$  much larger than data size

# Surrogate Model (continued)

- Explored point set:  $P_{exp} = \{\mathbf{x}_1, \dots, \mathbf{x}_p\}$ .

- Current responses:

$$V_{P_{exp}} = (f(x_1), \dots, f(x_p))^T$$

- Use  $\sum_j c_j \tilde{\phi}_j$  to approximate  $V_{P_{exp}}$ ,

$$\tilde{\phi}_j = (\phi_j(x_1), \dots, \phi_j(x_p))^T$$

- Two interesting questions:

- *Choice of the basis functions?*
- Estimation of the linear coefficients  $C_j$ ?

# Coefficient Inference

- Matching Pursuit Algorithm (Mallat and Zhang, 1993):

- Infer coefficients by minimizing  $\| V_{P_{\text{exp}}} - \sum_j c_j \tilde{\phi}_j \|$

- A greedy algorithm: at the  $j$ th iteration,

Let  $R^{(j-1)}$  be the current residual vector.

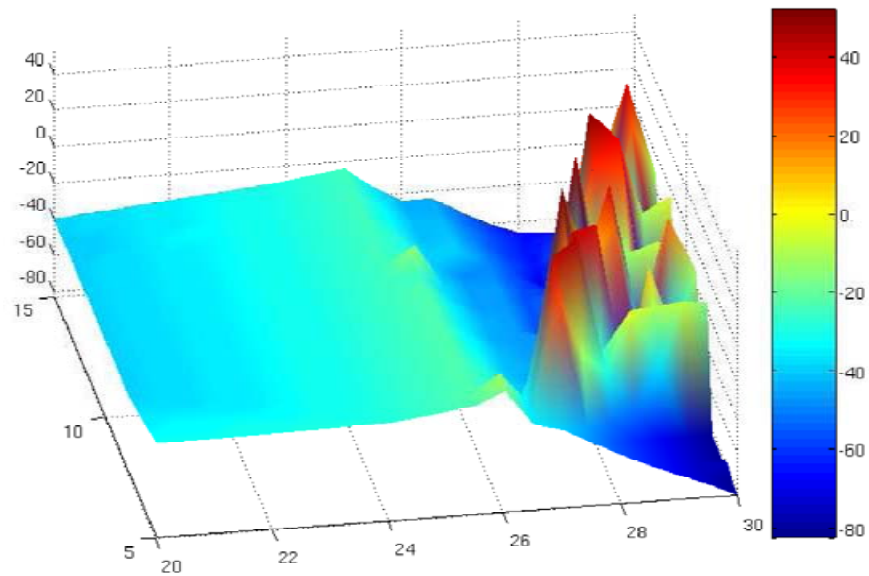
Selected a basis by  $\tilde{\phi}_{(j)} = \arg \max_i \langle R^{(j-1)}, \tilde{\phi}_i \rangle :$

$$c_{(j)} = c_{(j)} + \langle R^{(j-1)}, \tilde{\phi}_{(j)} \rangle,$$

$$R^{(j)} = R^{(j-1)} - \langle R^{(j-1)}, \tilde{\phi}_{(j)} \rangle \tilde{\phi}_{(j)}.$$

# Response Surface for Bistable Laser Diodes

- The true surface over a pre-specified grid:



- Search all positive Lyapunov exponents (PLE) (red area)
- PLE corresponds to chaotic light output.

# Gabor Functions

- Basis functions:

- $n$ -dimensional Gabor function

$$g(\mathbf{x}) = \exp\left(-\frac{\mathbf{x}^\top M \mathbf{x}}{2}\right) \exp(2\pi i A \mathbf{x}), \mathbf{x} = (x_1, \dots, x_n)^\top$$

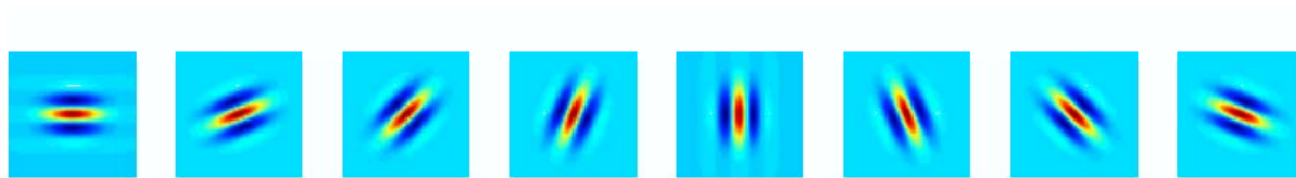
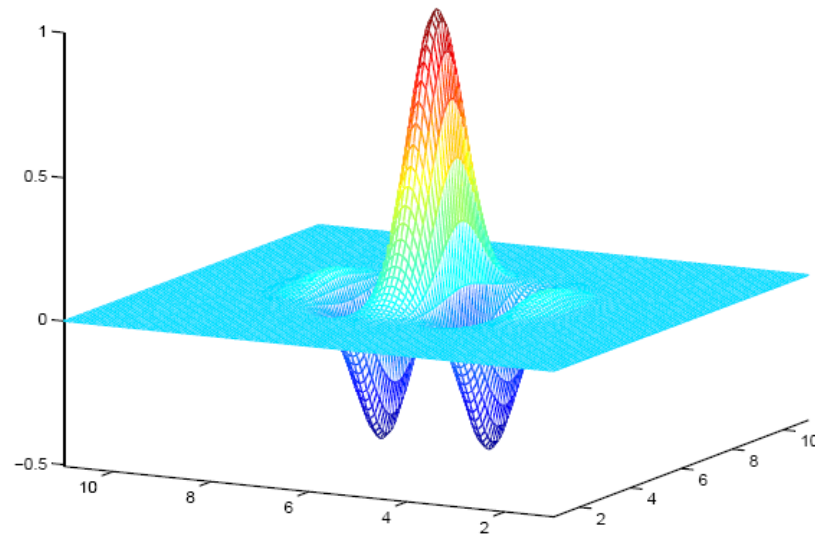
- Two-dimensional Gabor function, i.e.  $n = 2$

$$g(u, v) = \frac{1}{Z} \exp\left[-\frac{1}{2}(\sigma_u u^2 + \sigma_v v^2)\right] \cos\left[\frac{2\pi u}{\lambda} + \varphi\right],$$

$$u = u_0 + x_1 \cos \theta - x_2 \sin \theta$$

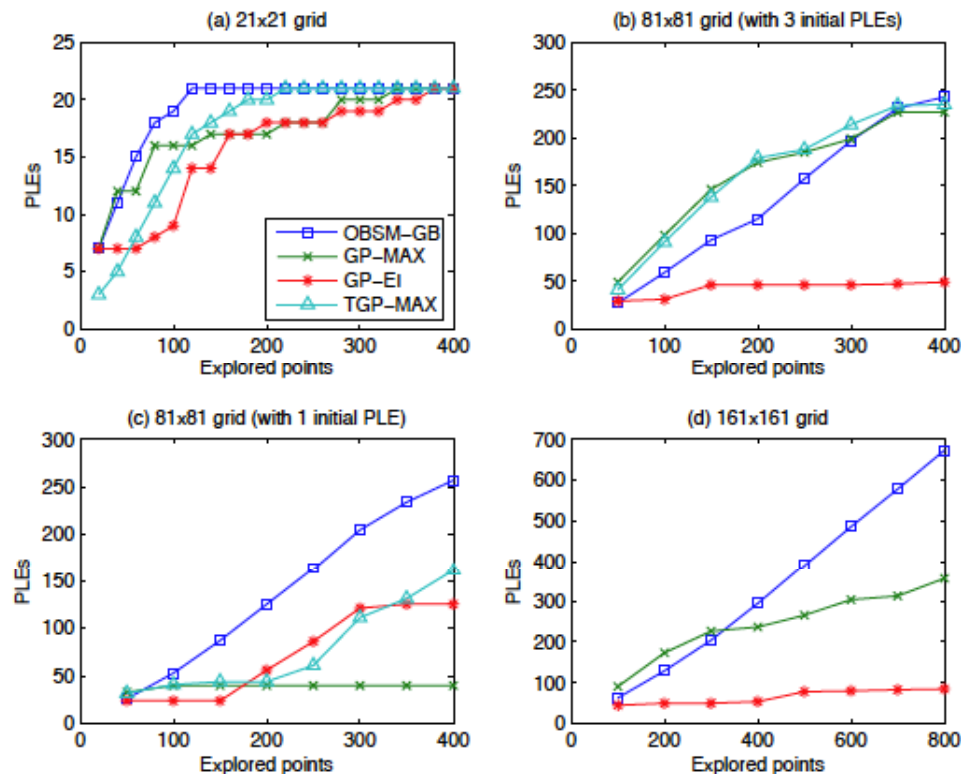
$$v = v_0 + x_1 \sin \theta - x_2 \cos \theta$$

# Plots of 2-D Gabor Function



# Overall Comparisons

Figure 6: Cumulative numbers of PLEs found by using different explored points in  $21 \times 21$ ,  $81 \times 81$ , and  $161 \times 161$  grids.



# Summary

- Computer experiments/simulations have become popular in engineering and science
- Kriging is the most common method for statistical meta-model building but is more limited for large or complex problems
- Alternatives to kriging are being sought:
  - Tweaking of kriging to achieve stability (regularized, hybrid, tapering, reduced rank)
  - Approximations with fast computations (OBSM, RIDW): but *lacking inferential capability*

# Covariance Matrix Tapering

- Covariance tapering (Kaufman et al., 2008)
  - ✓ Covariance matrix is “tapered” or multiplied element wise by a sparse matrix, to approximate the likelihood.
- Advantages:
  - ✓ Significant computational gains/stability.
  - ✓ Retain interpolating property.
  - ✓ Asymptotic convergence of the taper estimator.
- But:
  - ✓ The tapering function is isotropic: OK for spatial statistic problems, but not applicable to engineering problems.
  - ✓ The tapering radius needs to be determined.

# Rank Reduction

- Fixed rank kriging (Cressie-Johannesson, 2008)
  - ✓ A flexible family of non-stationary covariance function is defined by using a set of basis functions that are fixed in number (smaller than the data size  $n$ ).
- Advantage:
  - ✓ Reduce the computational cost of kriging to  $O(n)$ .
- But:
  - ✓ How to choose the appropriate basis functions.
  - ✓ Not an interpolator.

# Upper bound

- Upper bound

$$\kappa_2^p(\mathbf{R}(\boldsymbol{\theta}) + \lambda \mathbf{I}_n; \mathbf{X}) \leq \prod_{\substack{k=1 \\ \exp(-\theta_k)=1 \Leftrightarrow \theta_k=0}}^p \kappa_2^1(\mathbf{R}(0) + \lambda \mathbf{I}_{n_k}; D_k)$$

- The worst case of a correlation matrix

## Inverse Distance Weighting (IDW)

- Inverse Distance Weighting (Shepard, 1968):

$$\hat{y}(\mathbf{x}) = \frac{\sum_{k=1}^n w_k(\mathbf{x}) y_k}{\sum_{i=1}^n w_i(\mathbf{x})}.$$

- $w_i(\mathbf{x}) = 1/d(\mathbf{x}, \mathbf{x}_i)^2.$
- $d(\mathbf{x}, \mathbf{x}_i) = \left\{ \sum_{j=1}^p (x_j - x_{i,j})^2 \right\}^{1/2}.$
- Simple computation but poor prediction.

## Regression-Based Inverse Distance Weighting (RIDW)

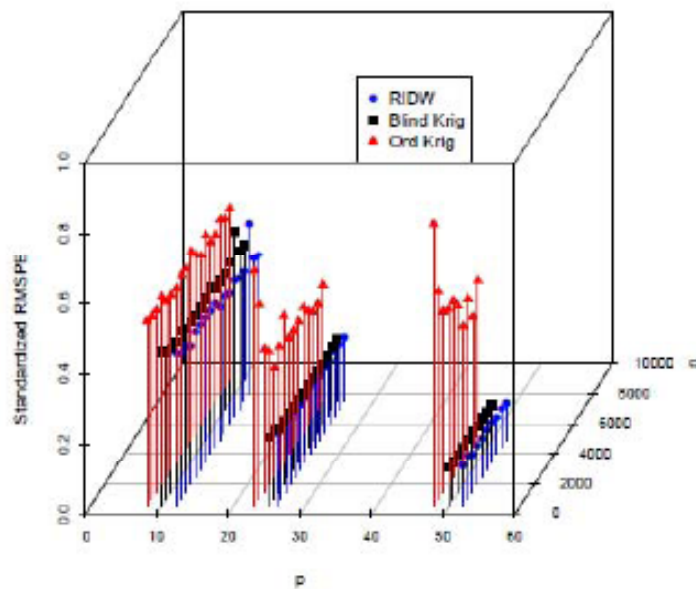
- Add regression part to IDW (Joseph and Kang, 2009):

$$\hat{y}(\mathbf{x}) = \mu(\mathbf{x}; \boldsymbol{\beta}) + \frac{\sum_{k=1}^n w_k(\mathbf{x}) e_k}{\sum_{i=1}^n w_i(\mathbf{x})}$$

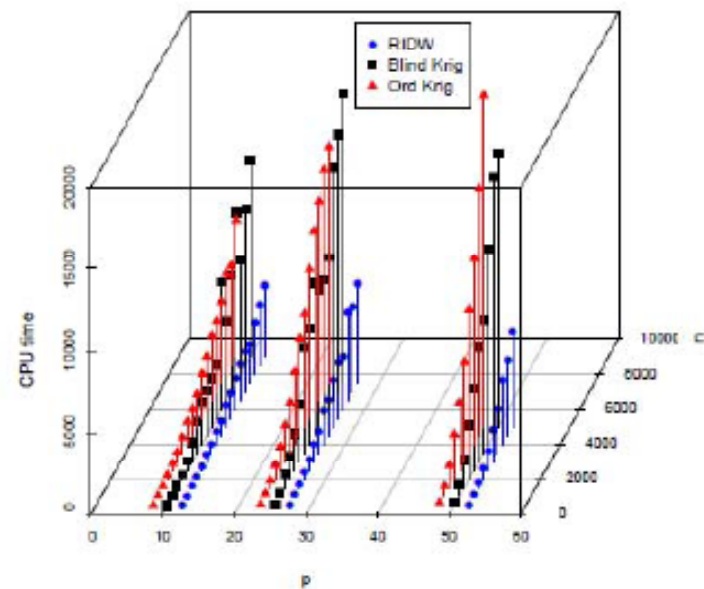
- $\mu(\mathbf{x}_k; \boldsymbol{\beta})$  = mean part; can be linear, nonlinear, nonparametric.
- $e_k = y_k - \mu(\mathbf{x}_k; \boldsymbol{\beta}) = y_k - \mu_k$ .
- $w_i(\mathbf{x}) = \frac{\exp\{-d^2(\mathbf{x}, \mathbf{x}_i)\}}{d^2(\mathbf{x}, \mathbf{x}_i)}$ . (faster convergence than IDW)
- $d(\mathbf{x}, \mathbf{x}_i) = \sqrt{\sum_{j=1}^p \theta_j (x_j - x_{i,j})^2}$ .

## Comparisons Between RIDW and Kriging

Standardized RMSPE



CPU time in simulation



# Lower bound on $\lambda$

- Lower bound

$$\lambda^* = \inf \left\{ \lambda \mid \prod_{k=1}^p (1 + n_k/\lambda) < \Delta \right\}$$

where

$$\epsilon = 2^{-52} \approx 2.22 \times 10^{-16}$$

Machine accuracy  
(or unit round-off)

$$\Delta = 1/(10\epsilon) = 4.5 \times 10^{14}$$