

# Thesis Defence

Simon Fraser University



Parameter Estimation and Uncertainty Quantification Applied to  
Advection-Diffusion Problems Arising in Atmospheric Source  
Inversion

Juan Gabriel García

April 10 2018

# Content

- 1 Overview
- 2 Setting up the Problem
  - Finding a Surrogate for  $F$
  - Optimizing the Surrogate Interpolation Capabilities
  - Reducing the Complexity of the Model
- 3 Introducing the Bayesian Framework
- 4 Decoding the Posterior

# Problem of Interest

Given a function

$$F : A \times \Theta \rightarrow \mathbb{R},$$

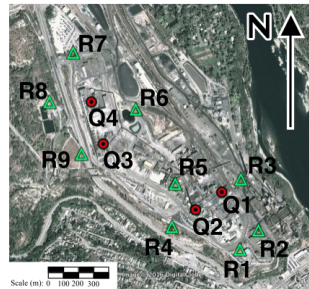
where  $A \subset \mathbb{R}^n$  and  $\Theta$  is a set of parameters. Given experimental (noisy) measures of  $F$  at known points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in A$ . How to infer the values of the parameters in  $\Theta$  and their uncertainties when  $F$  is computationally expensive?

# Case Study

Consider the model of pollutant transport for the concentration  $c(x, y, z, t)$  of a pollutant

$$\partial_t c + L(\theta)c = f$$

Goal: Estimate  $Q_i$  and  $\theta$  using measurements of deposition in  $R_j$ .



# Mathematical Model

$$\partial_t c(\mathbf{x}, t) + \nabla \cdot (\mathbf{u}(\mathbf{x}, t) c + \mathbf{S}(\mathbf{x}, t) \nabla c) = q(\mathbf{x}, t) \quad \text{on } \mathbb{R}^2 \times \mathbb{R}_{\geq 0} \times (0, T).$$

- $\mathbf{u}(\mathbf{x}, t) = (u_x(z, t), u_y(z, t), u_{set})$  (Wind velocity field)
- $\|(u_x, u_y)\|_2 \propto (z)^\gamma$ ,
- $\mathbf{S} = \text{diag}(s_x, s_y, s_z)$  (Eddy diffusion matrix),
- $s_z = f(L, z_{cut})$ ,
- $s_x = s_y = g(z_i, L)$ , with  $z_i$  Mixing layer height.

# Boundary Conditions

- Far-field boundary condition

$$c(\mathbf{x}, t) \rightarrow 0 \text{ as } \|\mathbf{x}\| \rightarrow \infty$$

- Robin boundary conditions at  $z = 0$

$$\left( u_{set} c + s_z \frac{\partial c}{\partial z} \right) \Big|_{z=0} = u_{dep} c \Big|_{z=0}$$

- To avoid inconsistencies in the Robin B.C. we define a cutoff length  $z_{cut}$ .

- Concentration and deposition are related via

$$w(x, y, T) = \int_0^T c(x, y, 0, t) u_{set} dt. \quad (1)$$

$$F(x_i, y_i, T) = \int_{R_i} w(x, y, T) dx dy \approx w(x_i, y_i, T) \Delta A,$$

- $\Delta A$  is the cross-sectional area of the dust-fall jar
- $T$  is taken to be one month.

- Numerical solution of the concentration  $c$  was obtained via a finite volume solver<sup>1</sup> using a 30x30 resolution grid on the domain.
- Solving for one set of parameters can take up to half an hour.

**Conclusion:** Finding the deposition  $F$  is computationally very expensive.

---

<sup>1</sup>Hosseini and Stockie, Computers and Fluids, 2017



# Roadmap

- 1 Find a surrogate for  $F$ ,
- 2 Locate optimal points to evaluate  $F$  so that surrogate is accurate,
- 3 See if it is possible to do dimensionality reduction,
- 4 Use the Bayesian framework to obtain the posterior distribution of parameters in the light of experimental data,
- 5 Perform inference on the posterior using numerical methods.

# Content

- 1 Overview
- 2 Setting up the Problem
  - Finding a Surrogate for  $F$
  - Optimizing the Surrogate Interpolation Capabilities
  - Reducing the Complexity of the Model
- 3 Introducing the Bayesian Framework
- 4 Decoding the Posterior

# Gaussian Process

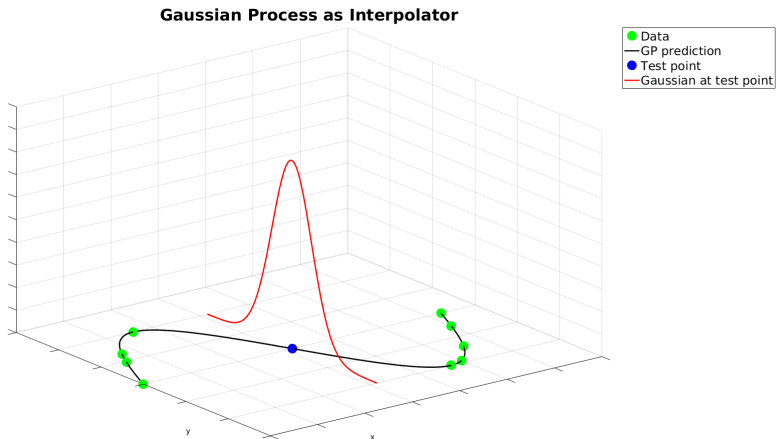
A Gaussian process (GP) is a collection of random variables  $\{g(x)\}_{x \in A}$ , for some set  $A$ , possibly uncountable, such that any finite subset of random variables  $\{g(x_k)\}_{k=1}^N \subset \{g(x)\}_{x \in A}$  for  $\{x_k\}_{k=1}^N \subset A$  are jointly Gaussian.

A GP is completely defined by its mean  $m(x)$  and covariance operator  $k(x, x')$ :

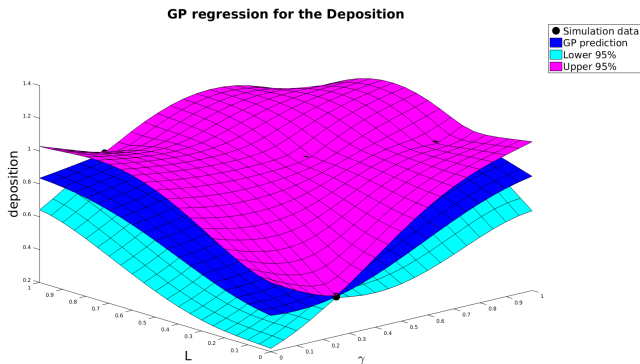
$$m(x) = \mathbb{E}(g(x)),$$

$$k(x, x') = \mathbb{E}((x - m(x))(x' - m(x'))).$$

# Gaussian Process as Interpolator

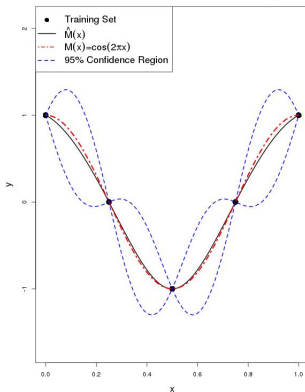


# Interpolation with Real Data

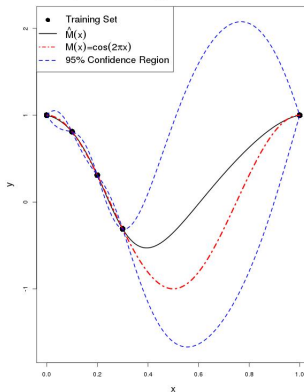


# Why Experimental Design?

Interpolation Using a Maximin Design



Interpolation Using an Arbitrary Partition

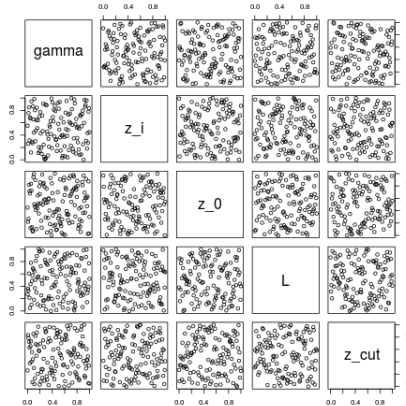


# Maximin Design

Given  $T \subset \mathbb{R}^n$  and a subset  $S$  of  $T$ , with finite (fixed) cardinality, say  $|S| = n$ . A maximin distance design  $S^o$  is a collection of points of  $T$  such that

$$\max_{(S \subset T, |S|=n)} \min_{(s, s' \in S)} \|s - s'\| = \min_{s, s' \in S^o} \|s - s'\| = \max!,$$

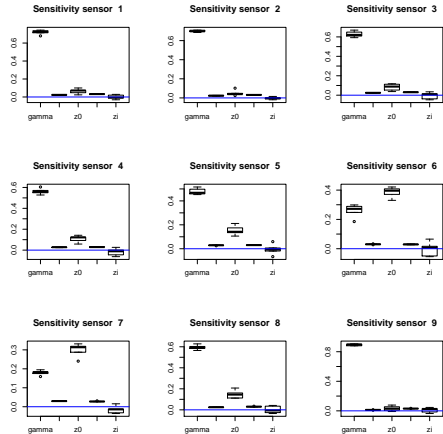
# Example of a Design





# Sensitivity Analysis

- $\gamma$ : Fitting parameter for the  $z$  dependence of the velocity.
- $z_0$ : Roughness length.
- $z_i$ : Mixing layer height.
- $L$ : Monin-Obukhov length.
- $z_{cut}$ : cutoff height.



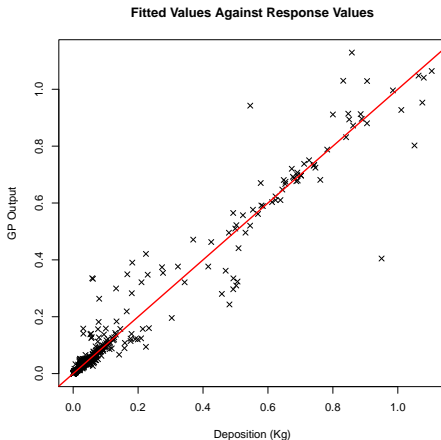
## What About the Sources?

The deposition behaves linearly with respect to the values of the 4 sources, hence

$$\begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_9 \end{bmatrix} = \mathcal{A}(\gamma, z_0, L) \begin{bmatrix} q_1 \\ \vdots \\ q_4 \end{bmatrix}$$

To evaluate the matrix  $\mathcal{A}(\gamma, z_0, L)$  we approximate its entries by Gaussian processes and create a surrogate  $A(\gamma, z_0, L)$ .

# Cross Validating the Surrogates



# Content

- 1 Overview
- 2 Setting up the Problem
  - Finding a Surrogate for  $F$
  - Optimizing the Surrogate Interpolation Capabilities
  - Reducing the Complexity of the Model
- 3 Introducing the Bayesian Framework
- 4 Decoding the Posterior

# Bayes' Rule

Given a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and two events  $A, B \in \mathcal{F}$ , with  $\mathbb{P}(B) \neq 0$ , we define the conditional probability of  $A$  given  $B$  by

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

and Bayes' formula

$$\mathbb{P}_{post}(A|B) \propto \mathbb{P}_{like}(B|A)\mathbb{P}_{prior}(A). \quad (2)$$

# Looking at the Stochastic Model

To account for the uncertainties in the interpolation and in the experimental measurements we propose the model

$$\begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_9 \end{bmatrix} = A(\gamma, z_0, L) \begin{bmatrix} q_1 \\ \vdots \\ q_4 \end{bmatrix} + \epsilon, \quad \text{where } \epsilon \sim \mathcal{N}(0, \lambda_\epsilon I_{9 \times 9})$$

# Probabilistic Model

Our goal is to estimate values of  $\omega := (\gamma, z_0, L)$  and  $q := (q_1, q_2, q_3, q_4)$  given measurements  $\vec{R}$ . Mathematically we want to estimate:

$$\mathbb{P}_{post}(\omega, q | \vec{R}) \propto \underbrace{\mathbb{P}_{like}(\vec{R} | \omega, q) \mathbb{P}_{prior}(\omega) \mathbb{P}_{prior}(q)}_{\text{Assuming } p \text{ and } q \text{ independent}}$$

We assume  $\omega \sim \text{Uniform}$  over the domain of definition of the parameters.

## What About $q$ ?

- $q_k > 0$  for  $k = 1, 2, 3, 4$ .
- If we trust the engineers the most likely value for  $q$  is the engineers estimate and the true value cannot be very far away from those estimates.

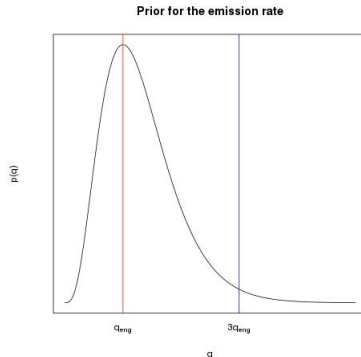
Source	Estimated Emission Rate [ton/yr]
$q_1$	35
$q_2$	80
$q_3$	5
$q_4$	5



## Choosing a Prior for $q$

A consistent assumption is  $q_k \sim Ga(\alpha_k, \beta_k)$ , for  $k = 1, 2, 3, 4$ , the following conditions that define  $\alpha_k$  and  $\beta_k$  for all  $k$  uniquely.

- $\beta_k(\alpha_k - 1) = q_{eng,k}$
- $qgamma(0.99, \alpha_k, \beta_k) = 3q_{eng,k}$



# Likelihood

Since  $\epsilon \sim \mathcal{N}(0, \lambda_\epsilon I_{9 \times 9})$  then

- $\mathbb{P}(\epsilon) \propto \exp\left(-\frac{\|\epsilon\|_2^2}{2\lambda_\epsilon^2}\right)$
- $\mathbb{P}_{like}(\vec{R}|\omega, q) \propto \exp\left(-\frac{1}{2\lambda_\epsilon^2}\|\vec{R} - A(\omega)q\|_2^2\right)$

# Calculating the Posterior

Putting everything together

- $\mathbb{P}_{like}(\vec{R}|\omega, q) \propto \exp\left(-\frac{1}{2\lambda_\epsilon^2}\|\vec{R} - A(\omega)q\|_2^2\right)$
- $\mathbb{P}_{prior}(\mathbf{q}) \propto \prod_{k=1}^k q_k^{\alpha_k-1} \exp(\beta_k q_k)$
- $\mathbb{P}_{prior}(\omega) \propto \mathbf{1}_{[0,0.6] \times [0,3] \times [-600,0]}$

Then

$$\mathbb{P}_{post}(\omega, q|\vec{R}) \propto \mathbb{P}_{like}(\vec{R}|\omega, q)\mathbb{P}_{prior}(\omega)\mathbb{P}_{prior}(q)$$

# Content

- 1 Overview
- 2 Setting up the Problem
  - Finding a Surrogate for  $F$
  - Optimizing the Surrogate Interpolation Capabilities
  - Reducing the Complexity of the Model
- 3 Introducing the Bayesian Framework
- 4 Decoding the Posterior

# Sampling From a Probability Distribution

---

## Algorithm 1 Adaptive Metropolis-Hastings Algorithm

---

```

1: Choose an initial point  $(\omega_1, \mathbf{q}_1)$  in the support of  $\mathbb{P}_{post}(\omega, \mathbf{q} | \vec{R})$ 
2:  $\beta = 0.05$ 
3: for  $j = 2 : N$  do
4:   if  $j \leq 14$  then
5:     Draw  $u$  from  $\mathcal{N}((\omega_j, \mathbf{q}_j), \frac{0.01}{7} I_{7 \times 7})$ .
6:   else
7:     estimate the empirical covariance matrix  $\Sigma_j$  based on the samples generated so far
8:     Draw  $u$  from  $(1 - \beta)\mathcal{N}((\omega_j, \mathbf{q}_j), \frac{(2.38)^2}{7} \Sigma_j) + \beta\mathcal{N}((\omega_j, \mathbf{q}_j), \frac{(0.1)^2}{7} I_{7 \times 7})$ .
9:   end if
10:  Propose  $(\tilde{\omega}_j, \tilde{\mathbf{q}}_j) \leftarrow (\omega_{j-1}, \mathbf{q}_{j-1}) + u$ .
11:  Compute  $\beta \leftarrow \min \left( 1, \frac{\mathbb{P}_{post}(\omega_j, \mathbf{q}_j | \mathbf{y})}{\mathbb{P}_{post}(\omega_{j-1}, \mathbf{q}_{j-1} | \mathbf{y})} \right)$ .
12:  Draw  $w \sim U([0, 1])$ .
13:  if  $w < \beta$  then
14:     $(\omega_j, \mathbf{q}_j) \leftarrow (\tilde{\omega}_j, \tilde{\mathbf{q}}_j)$       (Accept the move)
15:  else
16:     $(\omega_j, \mathbf{q}_j) = (\omega_{j-1}, \mathbf{q}_{j-1})$       (Reject the move)
17:  end if
18: end for

```

---

## Setting $\lambda_\epsilon$

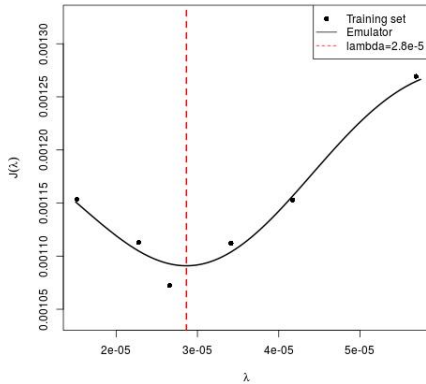
We define

$$J(\lambda_\epsilon) = \frac{1}{2} \int \left( \|A(\omega)q - \vec{R}\|_2 + \|q - q_{est}\|_2 \right) d\mathbb{P}_{post}^{\lambda_\epsilon},$$

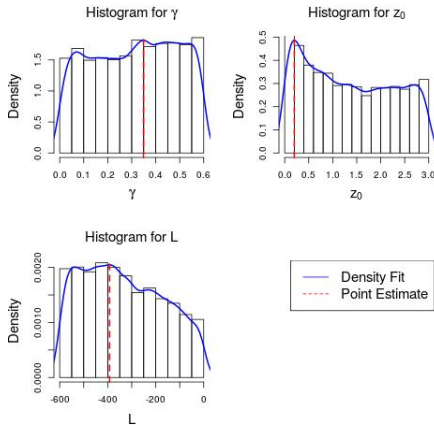
and choose a minimizer of  $J$

$$\hat{\lambda}_\epsilon = \operatorname{argmin} J(\lambda_\epsilon).$$

# Calculating $J$



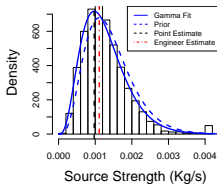
# Histograms for the Parameters



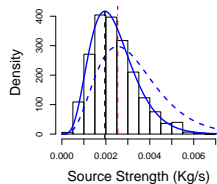


# Histograms for the Sources

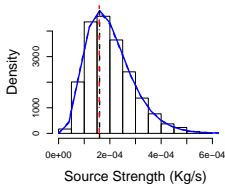
Histogram for Source 1



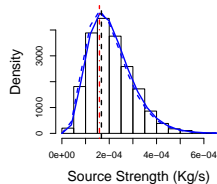
Histogram for Source 2



Histogram for Source 3



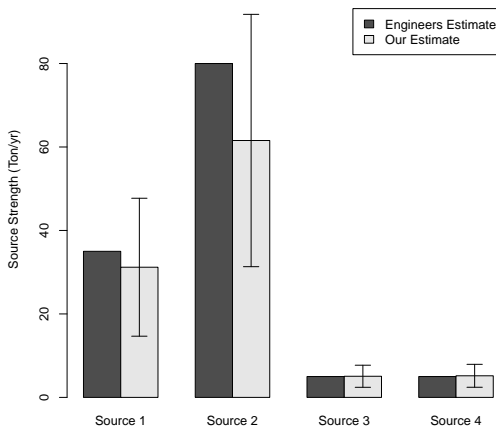
Histogram for Source 4



## Results for the Parameters

Parameter	Point Estimate	68% Confidence Interval
$\gamma$	0.3478	[0.1498, 0.5458]
$z_0$	0.0811	[0, 1.5781]
$L$	-379.45	[-195.86, -563.04]

## Results for the Sources



## Comparison with Related Work

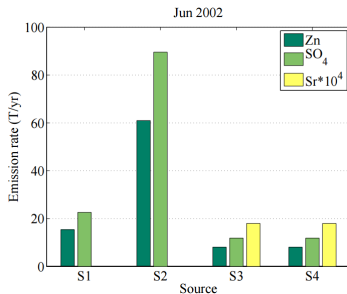


Figure: Lushi, E. and Stockie, J.M.  
Atmospheric Environment, 2010

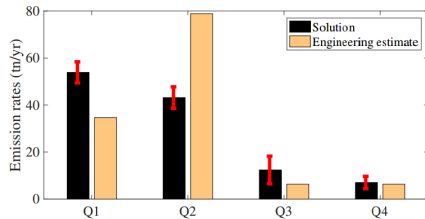


Figure: Hosseini and Stockie,  
Computers and Fluids, 2017

# Conclusions

- We have developed a method to cheaply estimate parameters in computationally expensive models.
- Instead of trial and error, we propose a methodology that allows estimating parameters in complex models using experimental data.
- Besides a point estimate we are able to obtain a confidence interval for it.
- Our results agree well with previous results.