# GAUSSIAN PROCESS REGRESSION

MATH-414: STOCHASTIC SIMULATION

Francesco Pettenon
Francesca Venturi

11th January 2023

# 1 INTRODUCTION

## 1.1 CONTEXT

This project finds its greatest motivation in the field of geoscience, where approximating unknown properties of the subsurface by means of Gaussian processes is a standard procedure. In particular, we know that scientists have at their disposal limited data collected through boreholes. The goal of the project is therefore to derive information about the entire permeability field of the area being studied by means of Gaussian Process Regression, that is, to make use of conditional Gaussian processes in order to reliably predict the unknown permeability field.

## 1.2 MOTIVATION

Gaussian process regression is a nonparametric method for regression that includes uncertainty quantification in its predictions and performs well on small datasets. It is based on Bayesian statistics, which assumes that the parameters describing the unknown field follow a known distribution (called the prior distribution) and that the observations also follow a certain distribution. Using Bayes' Rule, the prior distribution is updated based on the observations to obtain the posterior distribution, which incorporates information from both the prior and the observations. The predictive distribution is then obtained by integrating the posterior distribution and it is Gaussian as well. This method does not provide pointwise predictions, but rather a distribution from which single values can be derived using the mean and covariance matrix [1].

# 2 THEORETICAL BACKGROUND

This project focuses on the generalization of the above procedure to Gaussian processes and Gaussian random fields, known as Gaussian Process (GP) regression. The theoretical background is analyzed in detail in [1] and [2], which we have referred to throughout the entire project. For convenience, we list below the essential formulas and their notation. Let $\mathbf{f}(Z) \sim \mathcal{N}\left(\mathbf{m}(Z), \mathbf{K}(Z)\right)$ be a Gaussian random field, and denote with $Z$ any finite set of positions and $Y$ any new set of locations, whose value $\mathbf{f}(Y)$ is to predict. Furthermore assume that the we deal with noisy observations $\tilde{\mathbf{f}}(Z) = \mathbf{f}(Z) + \epsilon$ where $\epsilon \sim \mathcal{N}\left(\mathbf{0}, s^2\mathbf{I}_{N_Z}\right)$, then we know that the covariance of $\tilde{\mathbf{f}}(Z)$ is $\mathbf{K}(Z) + s^2\mathbf{I}$. Therefore, we have that $\mathbf{f}(Y)|\tilde{\mathbf{f}}(Z) \sim \mathcal{N}\left(\tilde{\mathbf{m}}(Y \mid Z), \mathbf{K}(Y|Z)\right)$, whose mean and covariance matrix are:

$$\tilde{\mathbf{m}}(Y \mid Z) := \mathbb{E}[\mathbf{f}(Y) \mid \tilde{\mathbf{f}}(Z)] = \mathbf{m}(Y) + \mathbf{K}(Y, Z)\left(\mathbf{K}(Z) + s^2\mathbf{I}\right)^{-1}\left(\tilde{\mathbf{f}}(Z) - \mathbf{m}(Z)\right) \tag{1}$$

$$\mathbf{K}(Y \mid Z) = \mathbf{K}(Y) - \mathbf{K}(Y, Z)\left(\mathbf{K}(Z) + s^2\mathbf{I}\right)^{-1}\mathbf{K}(Z, Y). \tag{2}$$

In our project we always assume the prior means $\mathbf{m}(Y) = 0$ and $\mathbf{m}(Z) = 0$.

## 2.1 COVARIANCE KERNELS

In the following study, we will make use of two covariance kernels, that is the tool for modeling covariance and the relationship between two variables. Both kernels we present are typically used to model stationary random processes. Recalling that $\|\cdot\|_2$ denotes the Euclidean norm of a vector, we introduce both kernels:
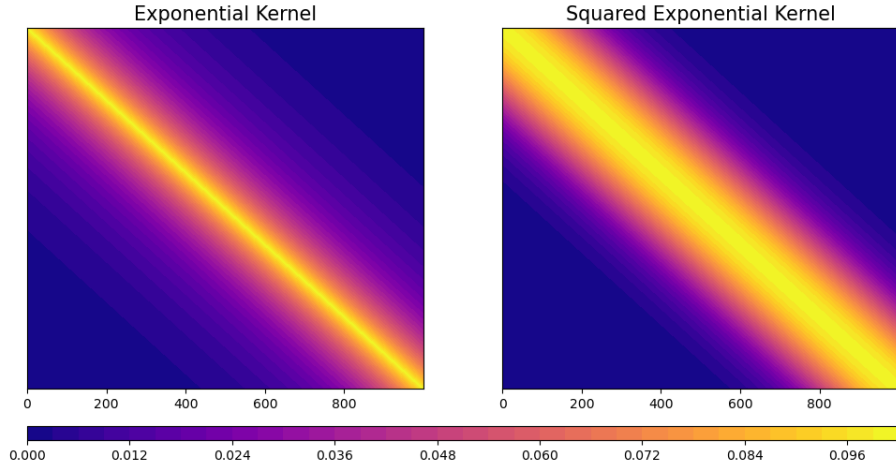
1. The Exponential (EXP) kernel given by

$$K_{exp}\left(\mathbf{x}, \mathbf{x}'\right) = \sigma^2 \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\ell}\right\}, \tag{3}$$

2. The Squared Exponential (SE) kernel given by

$$K_{se}\left(\mathbf{x}, \mathbf{x}'\right) = \sigma^2 \exp\left\{-\frac{1}{2}\sum_{i=1}^{k}\frac{(x_i - x_i')^2}{\ell_i^2}\right\} \tag{4}$$

In both of these expressions, we denote by $\sigma^2$ the variance and by $\ell = (\ell_1, ..., \ell_k)$ the correlation length or lengthscale. Although their formulation is similar, some considerations must be done.

One key difference between the two kernels is that the SE kernel gives more weight to points that are closer together by means of the square exponent, resulting in a smoother function. The EXP kernel, on the other hand, presents a linear exponent with respect to the distance between two points and this can lead to a rougher function. To be clearer, we plot two covariance matrices that are generated from the very same dataset, while they only differ with respect to the kernel employed: we can observe that the SE kernel has a larger correlation band around the diagonal, meaning that it actually gives more weight than the EXP to the points that are closer.



**FIGURE 1**
Visual difference between the EXP and the SE kernels 1D with fixed parameters $\sigma^2$ and $\ell$. These two matrices are taken as an example from Section (3) for a qualitative trend.

Another important difference arises in the case where the problem being studied is multimensional (see Section 4). On the one hand, the exponential kernel, being by construction isotropic, only has one value for $\ell$, whereas, instead, in the case of squared exponential kernel the correlation length becomes vector, thus admitting anisotropy. The isotropy of a kernel refers to whether the kernel is sensitive to the orientation of the data or not: as a consequence, the former is equally sensitive to the data in all directions, while the latter can capture directional patterns.

**HYPERPARAMETERS: $\ell$ AND $\sigma^2$**

**The Correlation Length $\ell$**: in both the EXP kernel and the SE kernel, the correlation length determines the degree of smoothness in the function being modeled. However, their interpretation is slightly different. In the exponential kernel, $\ell$ controls the rate at which the function being modeled drops off as the distance between two points increases. A larger correlation length means that the function will drop off more slowly, resulting in a smoother shape. A larger correlation length means that the function can vary over a larger length scale, resulting again in a smoother function.
**The Vertical Scale $\sigma^2$**: the vertical spread of a function is determined by its variance. The vertical scale $\sigma^2$ dictates the degree to which the predicted function deviates from the average prediction at any given point. A high variance implies greater flexibility, allowing the function to fluctuate more, while a low variance means that the function is restricted and will have less variations.

# 3   RECOVERING A SIMPLE FUNCTION

The preliminary stage of the study concerns recovering a scalar function. Indeed, suppose that the data we want to predict are generated by a simple function such as $f(x) = sin(x)$ and further assume that the $N_Z$ observations we have available are noise-free, i.e., that $y = sin(x)$. The goal of this first task is then to predict the entire Random Process, i.e., to generate $N_Y$ gaussian random variables, whose mean correctly approximates the sinusoidal function above.

## 3.1   GAUSSIAN PROCESS GENERATION

At the operational level, we proceed as follows.

- The domain of interest is the period of the function $sin(x)$ i.e., the interval $[0, 2\pi]$, from which we randomly extract $N_Z$ values $Z_i$ for $i = 1, ..., N_Z$, at which we evaluate the function $sin(Z_i)$, thus obtaining the Gaussian vector $\mathbf{f}(Z) \sim \mathcal{N}(\mathbf{m}(Z), \mathbf{K}(Z))$, where $Z = \{Z_i\}_{i=1}^{N_Z}$, $\mathbf{m}(Z) = sin(Z)$ and $\mathbf{K}(Z) = \mathbf{K}(Z, Z)$. Please note that, at this level, there is no difference between the two kernels and so, for now, we stick to a more general notation where the function $K(\mathbf{x}, \mathbf{x}')$ represents a generic kernel. We will look at the difference between $K_{exp}$ and $K_{se}$ when we analyze the results obtained (Section 3.2). Aiming to replicate the performance of a continuous function, we decide to divide the same interval into a sufficiently fine partition to hit the target, thus selecting a 1D uniform grid of $N_Y = 1000$ values, which we identify with $Y_i$ for $i = 1, ..., N_Y$.

- On the basis of what has already been presented in the section on theoretical analysis of Gaussian processes (Section 2), our intention is to study the conditional distribution $\mathbf{f}(Y)|\mathbf{f}(Z)$, of which we know the mean (1) and covariance (2). From these formulas, it is evident that, in order to generate the conditional Gaussian process, we need to compute the covariance matrices $\mathbf{K}(Z), \mathbf{K}(Y), \mathbf{K}(Y, Z)$. In this regard, therefore, we evaluate the function $K(\mathbf{x}, \mathbf{x}')$ for the vectors $Z$ and $Y$.

- Once the kernels have been evaluated, all that remains is to calculate the mean and variance of the conditional process. We note that although the observations are noise-free, $s^2 > 0$ is still considered. Indeed, we know that the computation of an inverse matrix can be numerically instable when the eigenvalues of the original matrix are close to zero. The term we add, therefore, allows each eigenvalue of that matrix to be increased by $s^2$, thus stabilizing the numerical procedure. Once the conditional mean and covariance are obtained, we observe that the covariance is nearly singular, so we proceed by generating the 3 separate Gaussian random processes using SVD, as described in [[2], Algorithm 3.1].
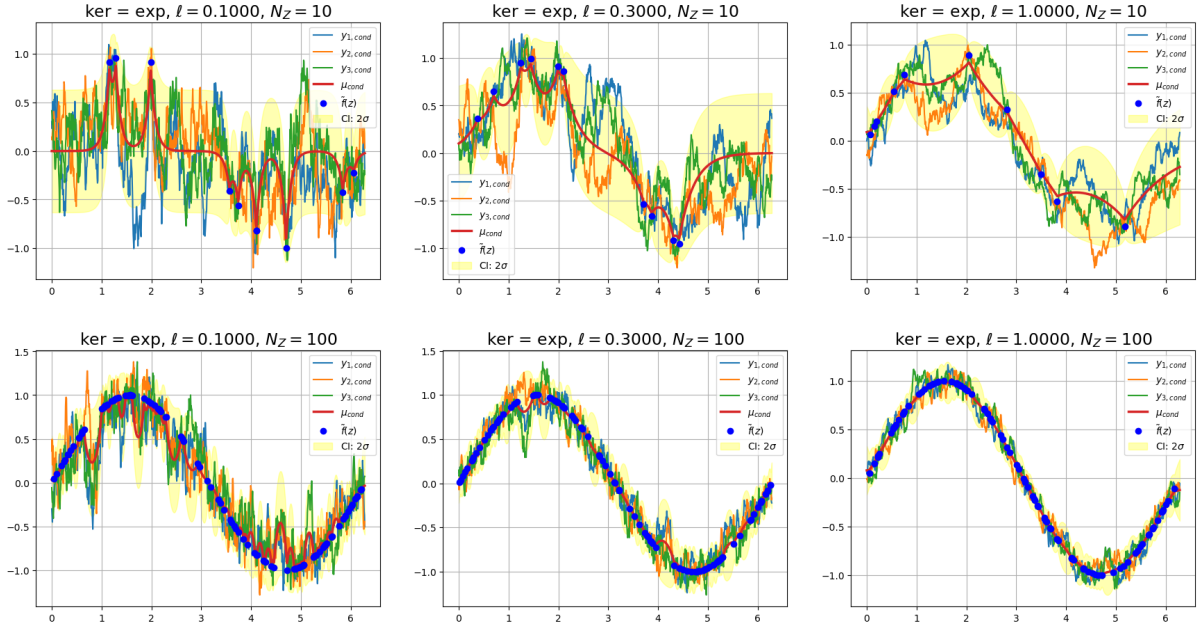
The whole procedure that we have just described is formally summarized in Algorithm 1, where the further step of optimizing the hyperparameters of the problem is performed, too.
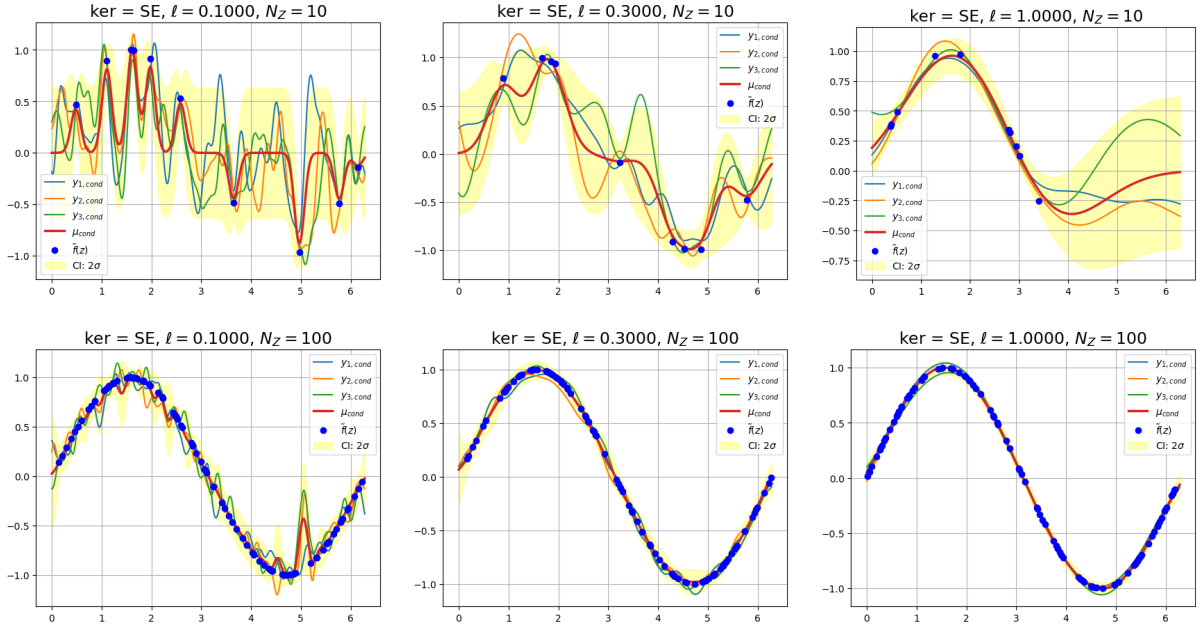
## 3.2   RESULTS

In this section we are concerned with analyzing the newly generated Gaussian processes and how they vary according to the amount of observations, i.e. $N_z$, and the choice of kernel.

Based on the plots generated in Figure (2-3), there are notable considerations:

- First, as is to be expected, the confidence interval decreases considerably as $N_Z$ increases: by adding information to the conditioning process (i.e., by increasing the vector of observations on which the model is trained) more meaningful information of the conditional distribution is obtained. This leads to a considerable decrease of the predictive uncertainty.

**FIGURE 2**
1D Gaussian Process: Exponential Kernel



**FIGURE 3**
1D Gaussian Process: Squared Exponential Kernel

- The most interesting observations, however, are about correlation length and the impact it has on posterior distribution. As noted in (2.1) the parameter $\ell$ characterizes the smoothness of the function being modeled. Indeed, from the plots, it is evident that as correlation length increases, the smoothness of the conditional mean (red line) increases. For example, for $\ell$=0.1, for the same distance between the variables, there is a poorer covariance: this implies that the posterior mean of the Gaussian processes does not deviate from the prior mean (zero in this case), except for the points at the $Z_i$ observations.

- As anticipated, the SE kernel assumes more correlation between points close to each other, so it

physiologically performs slightly better in our case.

- Finally, we observe that the confidence interval width does not cancel at $Z_i$ observations precisely because of the addition of the $s^2$ stabilizing term mentioned above.

## 3.3  OPTIMIZE HYPERPARAMETERS

In general, the specific parameters $\left\{\sigma^2, \ell, s^2\right\}$ that appear in the expression of the covariance kernels for GP regression can greatly impact the performance of the model. When no prior knowledge is available about their values, one approach to selecting them is to maximize the marginal likelihood

$$p_{\tilde{\mathbf{f}}(Z)}(v) = \int p_{\tilde{\mathbf{f}}(Z)|\mathbf{f}(Z)}(v \mid \mathbf{f})p_{\mathbf{f}(Z)}(\mathbf{f})\mathrm{d}\mathbf{f}, \text{ where } \begin{cases} p_{\tilde{\mathbf{f}}(Z)|\mathbf{f}(Z)}(\cdot \mid \mathbf{f}) = \mathcal{N}\left(\mathbf{f}, s^2\mathbf{I}\right) \\ p_{\mathbf{f}(Z)}(\cdot) = \mathcal{N}(\mathbf{0}, \mathbf{K}(Z)) \end{cases} \tag{5}$$

From this, one can compute the integral and obtain the marginal log-likelihood

$$L = \log p_{\tilde{\mathbf{f}}(Z)}(v) = -\frac{1}{2}v^T\left(\mathbf{K}(Z) + s^2\mathbf{I}\right)^{-1}v - \frac{1}{2}\log\left|\det\left(\mathbf{K}(Z) + s^2\mathbf{I}\right)\right| - \frac{N_Z}{2}\log 2\pi \tag{6}$$

As a result of the analytical computation of the marginal log-likelihood gradients (C), we obtain:

$$\frac{\partial L}{\partial s^2} = \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\operatorname{Tr}\left(\tilde{\mathbf{K}}^{-1}\right) \tag{7}$$

$$\frac{\partial L}{\partial \sigma^2} = \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\mathbf{K}_{\sigma^2}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\operatorname{Tr}\left(\tilde{\mathbf{K}}^{-1}\mathbf{K}_{\sigma^2}\right) \tag{8}$$

$$\frac{\partial L}{\partial \ell_i} = \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\mathbf{K}_{\ell_i}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\operatorname{Tr}\left(\tilde{\mathbf{K}}^{-1}\mathbf{K}_{\ell_i}\right) \tag{9}$$

Finally, we proceed by numerical optimization of the parameters by means of the gradient-based optimization algorithm `L-BFGS-B`, which yields the following values:
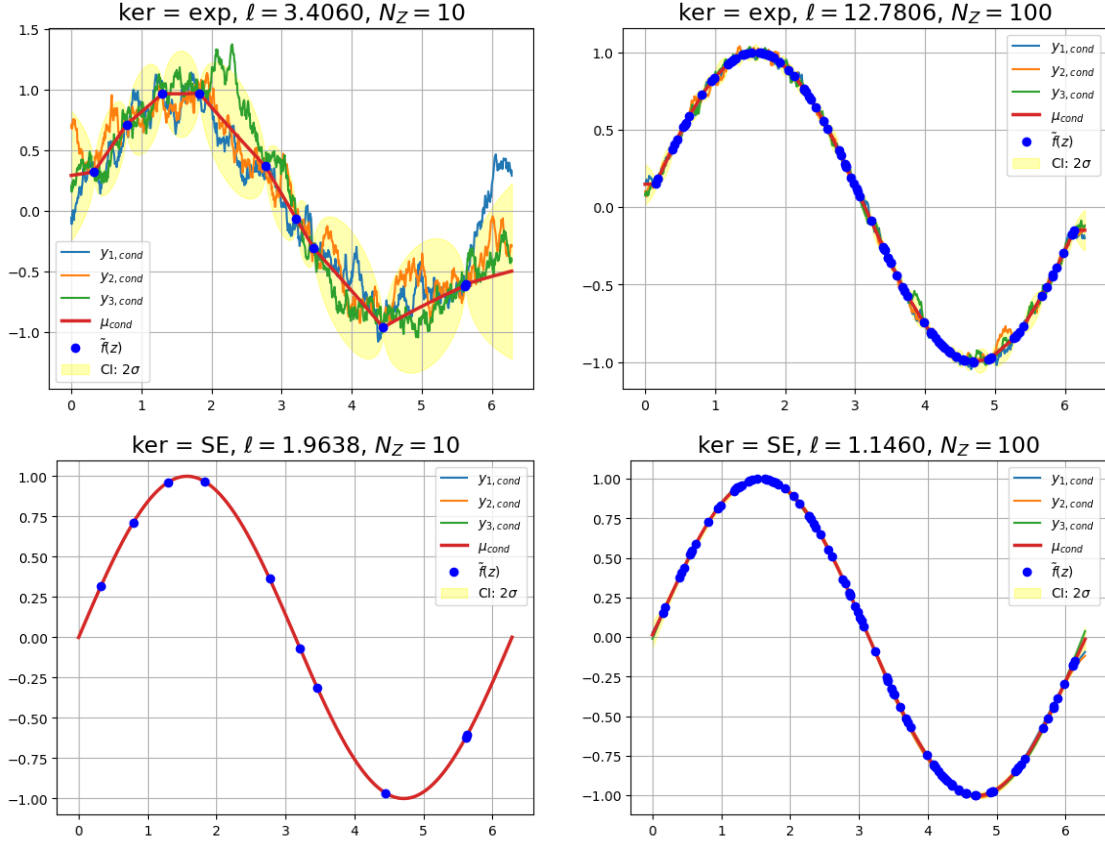
| $N_Z = 10$ | $\ell$ | $\sigma^2$ | $s^2$ |
|---|---|---|---|
| EXP | 3.406 | 0.411 | 0.000 |
| SE | 2.041 | 0.444 | 0.000 |

| $N_Z = 100$ | $\ell$ | $\sigma^2$ | $s^2$ |
|---|---|---|---|
| EXP | 12.758 | 0.182 | 0.000 |
| SE | 1.146 | 1.089 | 0.001 |

TABLE 1
Optimized hyperparameters

The most important feature is the optimized null value of $s^2$, from which we infer that the model learns that the observations on which it conditions are noise-free, so it recovers this feature even in the optimized case.

With the optimized values, we then resubmit the previous plots to visually capture the differences and improvements made by this procedure (Figure 4).

It thus becomes apparent how the absence of noise is reflected in the width of the confidence interval at the observations, which is now zero. Moreover, it is evident from the plots how the performance of the regression is extremely better, especially at a smaller dataset size. To conclude, the SE kernel, due to its smooth nature, will be more accurate in learning the smooth function $sin(x)$. On the other hand, the mean learnt by the EXP kernel tends to linearly interpolate between the training points.

**FIGURE 4**
Predictions and observations with optimized hyperparameters

The procedure we have just introduced generalizes 1D Gaussian Process Regression and is formalized in the following algorithm:

---
**Algorithm 1** 1D Gaussian Process Regression
---
1: Generate $N_Z$ Points Z uniformly distibuted in $[0, 2\pi]$
2: Generate 1000 Points Y in a 1D uniform grid of $[0, 2\pi]$
3: Given $f(x) = sin(x)$, evaluate the Points Z, obtaining $\mathbf{f}(Z)$
4: Optimize the parameters $\boldsymbol{\theta} = \left(\ell, \sigma^2, s^2\right)$ of the kernel w.r.t. marginal likelihood
5: Build the covariance matrices $\mathbf{K}(Y)$, $\mathbf{K}(Z)$, $\mathbf{K}(Y, Z)$
6: Compute the prediction $\mathbf{m}(Y|Z)$ and $\mathbf{K}(Y|Z)$ using (1) and (2)
7: Draw 3 indipendent realizations of the Gaussian process $\mathbf{f}(Y)|\mathbf{f}(Z)$

---

# 4 GP REGRESSION ON A PERMEABILITY FIELD

## 4.1 DATASET AND COMPUTATIONAL CONSIDERATIONS

In the second phase of the study, we apply the concepts learned in Section (3) to a geoscience problem in which we aspire to infer unknown properties of the subsurface permeability field based on a limited amount of available data. Specifically, we are provided with $60 \times 110$ observations subjected to noise $\epsilon \sim \mathcal{N}\left(\mathbf{0}, s^2\mathbf{I}_{N_Z}\right)$, representing the entire field. The goal of the study is to approximate this field by means of a much smaller number of data, which are used as training set for model generation. In particular, we will focus on two uniform rectilinear grids:

6

- The first dataset extracts an equispaced grid of $5 \times 5$ points, whose coordinates along $x$ are 10, 20, 30, 40, 50 and those along $y$ are 15, 35, 55, 75, 95

- The second dataset, on the other hand, consists of a grid of $11 \times 11$ observations, whose coordinates along $x$ are 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55 and along $y$ are 5, 15, 25, 35, 45, 55, 65, 75, 85, 95

The choice to analyze less data than the total amount is evidently dictated by the interest in generating a reliable predictive model that can perform properly once tested on a new sample, decreasing the cost of data collection and the computational complexity of the problem.

Before we begin the actual analysis of the problem, however, it is necessary to manipulate the data we have. In fact, the problem we are going to analyze is in two dimensions: we need to refer to a one-dimensional vector of observations, which allows us to apply the techniques analyzed in the Section (3). All the matrices we will deal with are inspected in *Row-Major Order*: elements of a multi-dimensional array are arranged sequentially row by row, which means filling all the index of first row and then moving on to the next row. What is thus obtained are three vectors: the entire permeability field `true_perm` consisting of 6600 real observations, the first dataset `Z1` of 25 points and the second dataset `Z2` of 121 points.

## 4.2 OPTIMIZE HYPERPARAMETERS

In order to optimize the hyperparameters, we use the same method as before, i.e., marginal log-likelihood maximization, obtaining the following results. Notice that, in this case, the SE kernel has two correlation length values, thus it is able to pick up any directional trend in the permeability field.

| Z1 | $\sigma^2$ | $s^2$ | $\ell_1$ | $\ell_2$ | Z2 | $\sigma^2$ | $s^2$ | $\ell_1$ | $\ell_2$ |
|------|--------|-------|--------|--------|------|--------|-------|--------|--------|
| EXP | 12.395 | 1e-10 | 64.590 | - | EXP | 11.428 | 0.035 | 40.395 | - |
| SE | 9.047 | 0.589 | 16.330 | 21.891 | SE | 7.005 | 1.120 | 11.720 | 18.266 |

TABLE 2
Optimal hyperparameters $\theta = \left\{ \sigma^2, s^2, \ell, \right\}$

## 4.3 PREDICT THE GAUSSIAN PROCESS

Once the parameters have been optimized, it is possible to obtain the matrices $\mathbf{K}(Y)$, $\mathbf{K}(Z)$, $\mathbf{K}(Y, Z)$ for both kernels and both domains. Finally, by exploiting the formulas for the distribution of conditional Gaussian random variables, the mean $\tilde{\mathbf{m}}(Y \mid Z)$ and covariance matrix $\mathbf{K}(Y|Z)$ can be computed. So, for each kernel and for each dataset, the algorithm is:
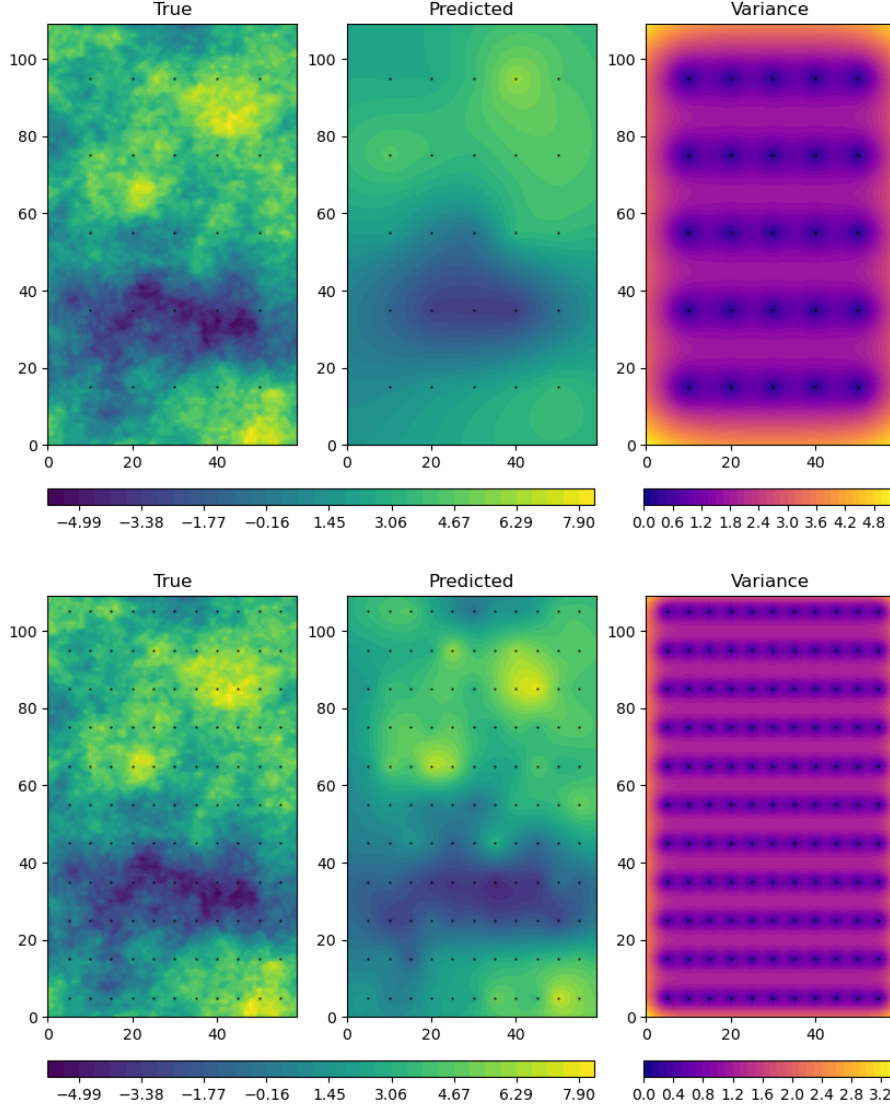
---
**Algorithm 2** 2D Gaussian Process Regression
---
1: Create 2D-meshgrid with Points $Z$
2: Extract the Points Z in `true_perm`, obtaining $\mathbf{f}(Z)$
3: Create the link between the physical points (dim = $(110 \times 60)$) and the Gaussian Vector $Y$ (dim = 6600):

$$Y[n] \iff \texttt{Points}\,[n/60]\,[n\%60] \quad \text{and} \quad \texttt{Points}\,[i]\,[j] \iff Y[i \cdot 60 + j]$$

4: Optimize the parameters $\boldsymbol{\theta} = \left(\boldsymbol{\ell}, \sigma^2, s^2\right)$ of the kernel w.r.t. marginal likelihood
5: Build the covariance matrices $\mathbf{K}(Y)$, $\mathbf{K}(Z)$, $\mathbf{K}(Y, Z)$
6: Compute the prediction $\tilde{\mathbf{m}}(Y \mid Z)$ and $\mathbf{K}(Y|Z)$ using (1) and (2)
7: **return** $\tilde{\mathbf{m}}(Y \mid Z)$, $\mathbf{K}(Y|Z)$
---

In Figure (5 - 6) we plot, for each kernel, the mean and variance of the approximate permeability field at both datasets of observations.
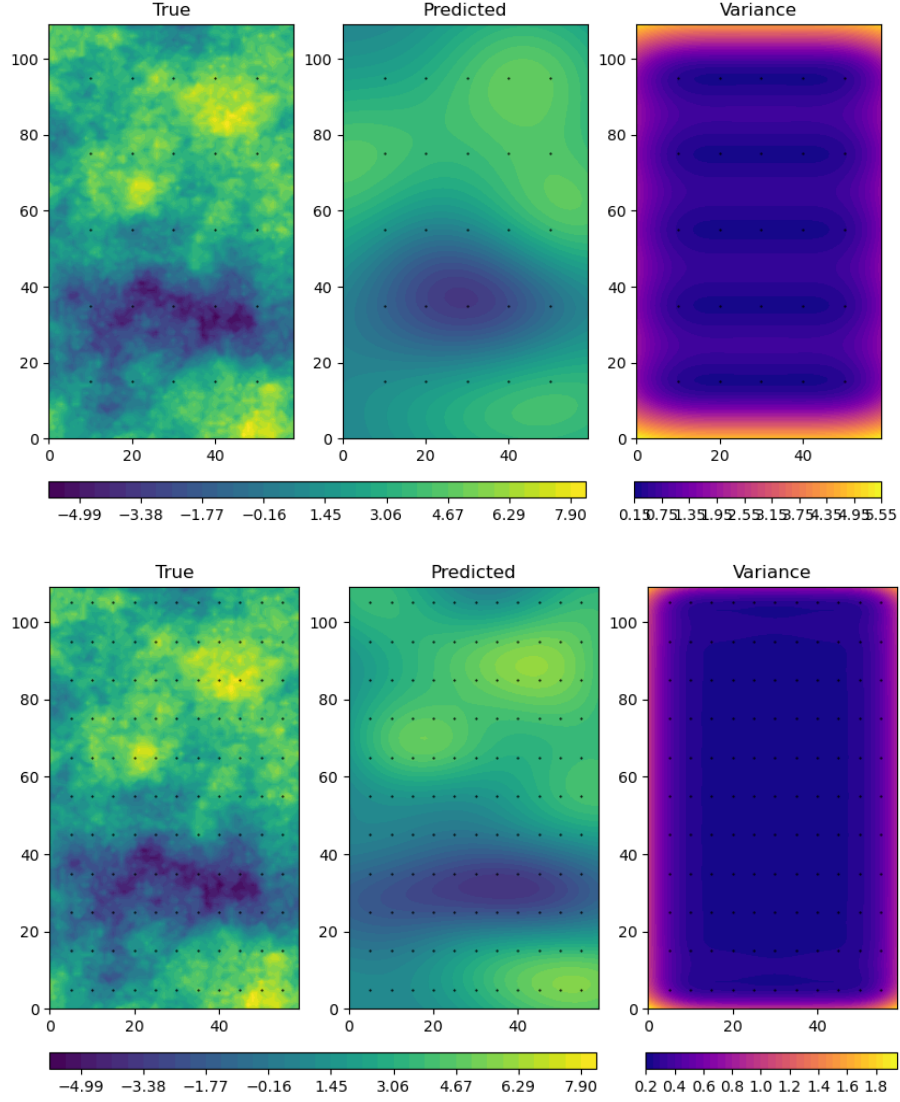


**FIGURE 5**
2D Gaussian Process: Exponential Kernel

Although plots in two dimensions are inherently less intuitive, some noteworthy behaviors can be detected.

- First, it is observed that by increasing the dimension of the training set, the variance of the observations falls, which means that the confidence interval of each prediction narrows. Intuitively, as the number of available data increases, the performance of the model improves, decreasing its uncertainty.

- In the plots of the variances we observe that, at the interpolation points, the variance decreases drastically, reaching $s^2$ values in each of them (except for errors of the order of $10^{-4}$): this is the behavior we expect in the presence of observations subject to noise. This behavior is more pronounced in the case of exponential kernels.

- In addition, it can be seen that the variance increases at the edges of the domain: this is typical behavior of Gaussian process regression.

- Finally, the exponential squared kernel causes the mean of the predictions to be a much smoother function: in fact, the contour lines are more regular than those for the exponential kernel.



**FIGURE 6**
2D Gaussian Process: Exponential Kernel

## 4.4 GENERATE PROCESSES: CIRCULANT EMBEDDING

To generate some realization of the 2D conditional Gaussian random fields $\mathbf{f}(Y)|\tilde{\mathbf{f}}(Z) \sim \mathcal{N}(\tilde{\mathbf{m}}(Y \mid Z), \mathbf{K}(Y|Z))$, there are several methods that can be applied. The standard methods involve factorizing the covariance matrix $\mathbf{K}(Y|Z)$ by means of Cholesky's method or SVD: this approach, however, may be computationally heavy when this matrix is large. For this reason we decide to resort to Circulant Embedding, that is, the generation of a circular matrix from the covariance matrix $\mathbf{K}(Y|Z)$, which enjoys convenient algebraic properties. It is in fact easily factorized by means of the FFT since its eigenvalues and eigenvectors are computable by known formulas. Once the circular matrix $\mathbf{P}$ has been generated, the information needed to generate a Gaussian process can be derived from it ([3], Chapter 3.3.2). However, the construction of the matrix $\mathbf{P}$ in 2D is not as straightforward as in the one-dimensional case: the notion of homogeneity and isotropy differ as soon as we move to dimensions greater than 1. In light of this difference, we find that the covariance matrix $\mathbf{K}(Y|Z)$ of the conditional process is not homogeneous

9

for both kernels: it is therefore not possible to proceed by deriving the circular matrix directly from it. Consequently, we choose to apply circulant embedding to the matrix $\mathbf{K}(Y)$ to obtain a vector of Gaussian random variables, using Agorithm 3.

---

**Algorithm 3** Circulant Embedding 2D

---

1: Build $\mathbf{P}$ starting from $\mathbf{K}(Y)$
2: Calculate $\mathbf{W} = \text{fft}(\mathbf{P})$, the 2-dimensional FFT of $\mathbf{P}$
3: Check that all elements of $\mathbf{W}$ are positive
4: Generate matrix $\mathbf{X}$ with size of $\mathbf{W}$ containing i.i.d. normal variables
5: Calculate $\tilde{\mathbf{Z}} = \text{ifft}(\sqrt{\mathbf{W}} \odot \mathbf{X})$, the 2-dimensional inverse FFT of $\sqrt{\mathbf{W}} \odot \mathbf{X}$
6: Generate $\mathbf{f}(Y) = Re(\tilde{\mathbf{Z}})|_Y + Im(\tilde{\mathbf{Z}})|_Y$
7: **return** $\mathbf{f}(Y)$

---

We observe that, within the 6600 values just generated, there are also the sampling points, whose true value is known. In light of this, "a posteriori" conditioning can be applied to the observations, using Algorithm 4.
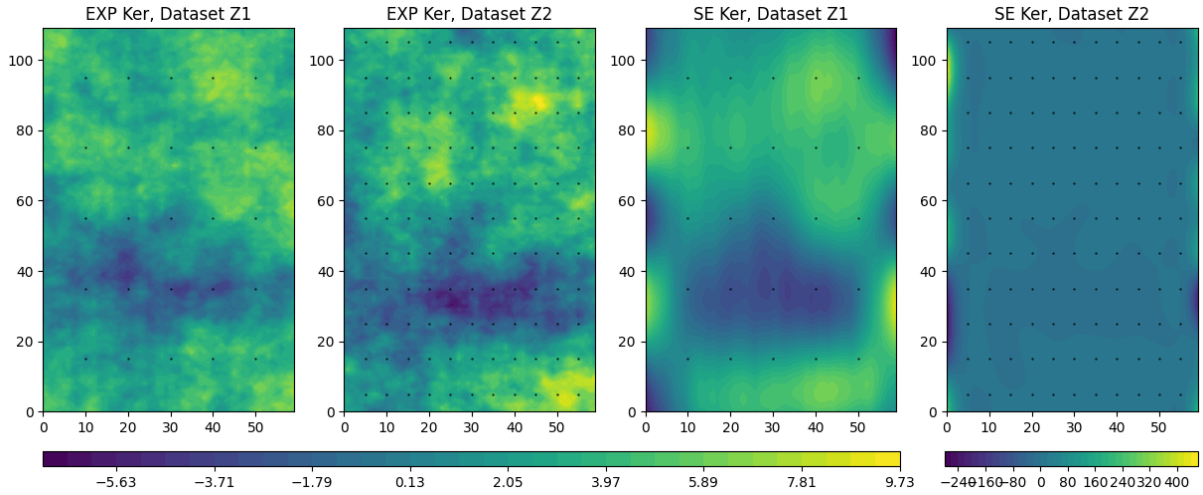
---

**Algorithm 4** Generation from conditional Gaussian distribution

---

1: Generate $\mathbf{f}(Y) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(Y))$ using 2D Circulant Embedding
2: Extract indexes of points $Z$ and evaluate $\mathbf{f_{new}}(Z) = \mathbf{f}(Y)|_Z$
3: Generate $\mathbf{f}(Y)|\mathbf{f}(Z) = \mathbf{f}(Y) + \mathbf{K}(Y, Z)\mathbf{K}(Z)^{-1}(\mathbf{f_{new}}(Z) - \mathbf{f}(Z))$
4: **return** $\mathbf{f}(Y)|\mathbf{f}(Z)$

---

As a result, we obtain the conditional Gaussian processes $\mathbf{f}(Y)|\tilde{\mathbf{f}}(Z) \sim \mathcal{N}(\tilde{\mathbf{m}}(Y|Z), \mathbf{K}(Y|Z))$, whose visual representation follows.



**FIGURE 7**
Gaussian Processed obtained via Circular Embedding technique

Observe that this procedure is extraordinarily unstable for the SE kernel when we deal with the second set of training points Z2.

## 4.5 MONTE CARLO ESTIMATOR

At this stage, our aim is to determine whether the permeability at certain coordinates exceeds a pre-determined threshold. Specifically, we are examining four critical locations with coordinates $x = $

$35, 40, 45, 50$ and $y = 85$ (and corresponding indexes $5135, 5140, 5145, 5150$). We want to estimate $\mathbb{P}(\max_{i=1,\dots,4} f(x_i, y_i) \geq 8) = \mathbb{E}(\psi(\mathbf{x}, \mathbf{y}))$, where $\psi(\mathbf{x}, \mathbf{y}) = \mathbf{1}_{\max_{i=1,\dots,4} f(x_i,y_i) \geq 8}$.

**Theorem 1** *Let $\mathbf{X}$ be a Gaussian vector of size $n$ with mean $\boldsymbol{\mu}$ and covariance matrix $\Sigma$. Then the linear transformation $\mathbf{Y} = A\mathbf{X} + \mathbf{b}$ will transform $\mathbf{X}$ into a Gaussian vector $\mathbf{Y}$ with mean $A\boldsymbol{\mu} + \mathbf{b}$ and covariance matrix $A\Sigma A^T$.*

The result of this theorem is that generating a Gaussian process limited to the four points of interest will always produce a Gaussian process whose mean and covariance are obtained evaluating the original one in those points. Therefore, we can use the Monte Carlo method (CMC) to generate a large number of processes and evaluate the probability with significantly lower computational cost. Using this method, we can estimate the probability and the $1 - \alpha$ confidence interval.

Given a tolerance for the range of the confidence interval, the number $N$ of the simulations to perform is chosen with a two stages Monte Carlo method, whose algorithm is described in [[2], Algorithm 5.1].

The results obtained can be found in Table (3):

| z2 | $\bar{N}$ | $\hat{\mu}_{\bar{N}}$ | $\hat{\sigma}_{\bar{N}}$ | tol | z2 | $\bar{N}$ | $\hat{\mu}_{\bar{N}}$ | $\hat{\sigma}_{\bar{N}}$ | tol |
|---|---|---|---|---|---|---|---|---|---|
| EXP | 9360182 | 0.024976 | 0.15605 | 1e-4 | EXP | 26223987 | 0.073693 | 0.26127 | 1e-4 |
| SE | 2633285 | 6.7596e-5 | 0.00822 | 1e-5 | SE | 768284 | 2.6032e-5 | 0.00161 | 1e-5 |

**TABLE 3**
Estimation with Two stages MC for datasets `z1` and `z2`

We notice that the simulations carried out using the EXP kernel tend to produce higher results than those using the SE kernel. This is to be expected, as the mean values at specific points with the EXP kernel are larger than those with the SE kernel, even though they are all less than $8$. Additionally, the variance is also higher for the EXP kernel, which means there is a higher probability of reaching the value of 8 due to a larger deviation from the mean. In contrast, this is a less common occurrence when using the SE kernel.

## 4.6 VARIANCE REDUCTION

One of the major drawbacks of the Crude Monte Carlo (CMC) method is its slow convergence rate of $O(\frac{1}{\sqrt{N}})$, resulting in a lengthy process to obtain precise results. One method to decrease the variance of Monte Carlo and consequently accelerating convergence is the use of Antithetic Variables (AV). AV involves using a probability distribution that is symmetrical about the origin to generate the random negatively correlated pairs of samples used in the simulation. This technique helps to reduce the approximation error and achieve more accurate results in a shorter amount of time.

**Theorem 2** *Assume that the random variable $Z$ has the expression $Z = \psi(X)$, with $X = (X_1, \dots, X_d)$ a random vector with independent components, such that*

- *$X$ has a symmetric distribution around its mean, i.e. $2\mathbb{E}[X] - X \sim X$*

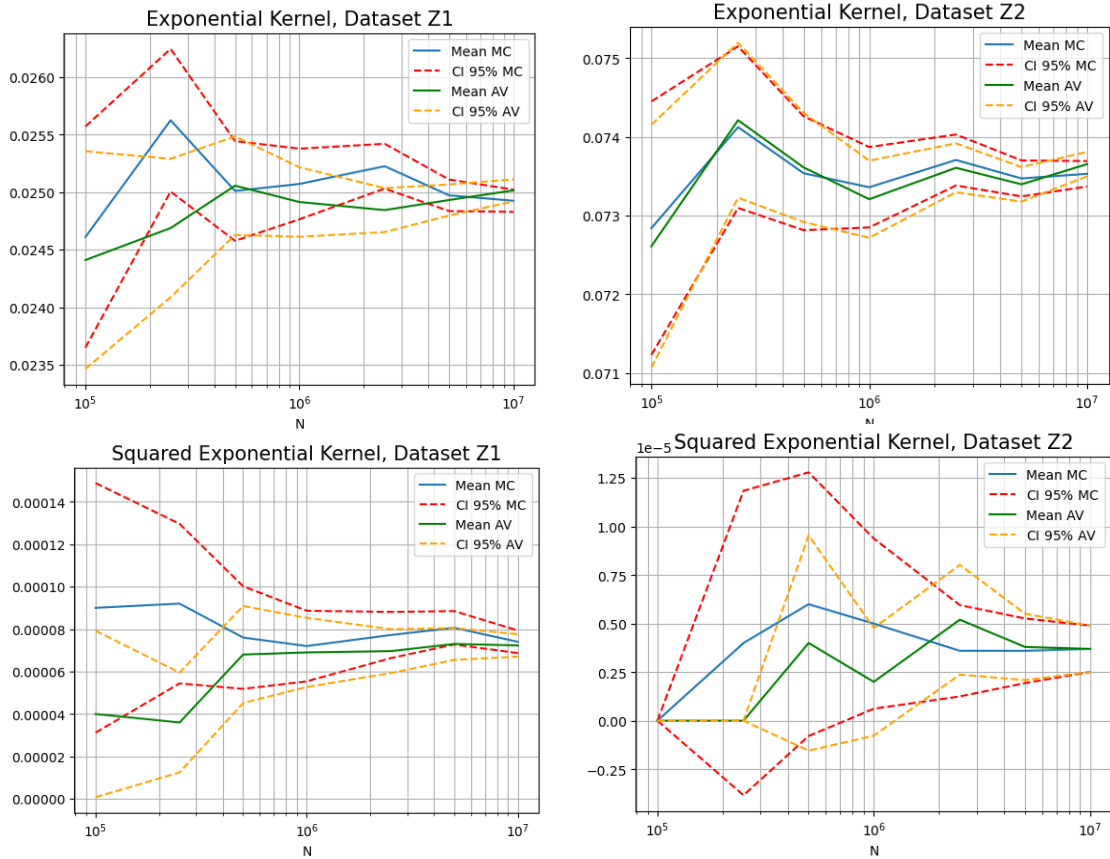- *$\psi$ is a monotone function in each of its arguments.*

*Then $Z = \psi(X)$ and $Z_a = \psi(2\mathbb{E}[X] - X)$ satisfy $\mathbb{E}[Z] = \mathbb{E}[Z_a]$ and $\text{Cov}(Z, Z_a) < 0$.*

It can be easily confirmed that the conditions required for this case are satisfied. Specifically, $\psi(\mathbf{X})$ is a characteristic function and so monotone with respect to each variable in the random vector $\mathbf{X} = (X_1, \dots, X_4)$. Additionally, since $\mathbf{X}$ is a Gaussian random vector, its distribution is symmetric around its mean $\boldsymbol{\mu}(X)$. A Monte Carlo approximation with antithetic variables can be constructed using [[2], Algorithm 6.1]. We generate different estimation of CMC and AV increasing N. In Table 4, we show the confidence interval for the last observations ($N = 1e8$).

| Z1 | $\hat{\mu}_{CMC}$ | $\hat{\mu}_{AV}$ | Z2 | $\hat{\mu}_{CMC}$ | $\hat{\mu}_{AV}$ |
|---|---|---|---|---|---|
| EXP | $0.0249 \pm 9.662e\text{-}5$ | $0.0250 \pm 9.560e\text{-}5$ | EXP | $0.0737 \pm 1.6193e\text{-}4$ | $0.0737 \pm 1.5539e\text{-}4$ |
| SE | $7.4e\text{-}5 \pm 5.331e\text{-}6$ | $7.23e\text{-}5 \pm 5.269e\text{-}6$ | SE | $3.7e\text{-}6 \pm 1.1922e\text{-}6$ | $3.7e\text{-}6 \pm 1.1922e\text{-}6$ |

TABLE 4

Confidence interval of CMC and AV for datasets `Z1` and `Z2`

In Figure (8), we plot the expected values predicted and the confidence intervals with respect to $N$, both for CMC and AV. The use of the antithetic variables method slightly decreases the confidence interval for the expected value compared to the crude Monte Carlo method, as demonstrated by our assumptions. However, when using the SE kernel with small sample sizes N, the estimate of the expected value becomes highly variable and may even be zero. To obtain a more accurate estimate, a larger sample size is required, which increases computational cost. The results show that the AV method is more effective when using the EXP kernel, rather than for modeling less probable events like the SE kernel.



FIGURE 8

Monte Carlo estimations and Antithetic Variables to reduce variance

## 5 CONCLUSION

The use of Gaussian Process Regression in this project has allowed us to successfully reconstruct a one-dimensional function (Section 3) and estimate the permeability field through two-dimensional observations in a geoscience problem (Section 4). However, we have observed that the choice of the kernel can significantly impact the performance of the process. This aligns with the principles of the Bayesian approach, in which the prior distribution carries significant weight. We showed the utility of Gaussian Process Regression in geoscience applications, but also emphasized the need to carefully consider the kernel selection in order to achieve optimal results. The implementation of the algorithms to carry out the project can be found in https://github.com/francpp/Gaussian-Process-Regression.

# A APPENDIX: SIMULATIONS

## A.1 TASK1: RECOVERING A SIMPLE FUNCTION

```python
# define initial parameters
NNz = [10, 100]
Ny = 1000
var = 0.1
s_2 = 0.01
ll = [0.1, 0.3, 1.0]
f_tilde = lambda z: np.sin(z)

# do simulations for all parameters
for ker in range(2):
    if ker == 0:
        kernel = lambda x1, x2, l, var: kernel_exp(x1,x2,l,var)
        type_ker = 'ker_exp'
    elif ker == 1:
        kernel = lambda x1, x2, l, var: kernel_SE(x1,x2,l,var)
        type_ker = 'ker_SE'

    for l in ll:
        for Nz in NNz:

            #create domain
            Z, f_tilde_Z, Y = create_Domain1(Nz, f_tilde, Ny)
            #evaluate kernels
            K_Z, K_Y, K_ZY = evaluate_kernels(Z, Y, kernel, l, var)
            #computes conditional mean and variance of the process
            mu_YgivenZ, K_YgivenZ = prediction(f_tilde_Z, K_Z, K_Y, K_ZY, s_2)
            #samples again
            y_cond=g_process(mu_YgivenZ, K_YgivenZ, n_process = 3)

# Minimize THETA of the log-likelihood

Z0, f_tilde_Z0, Y = create_Domain1(NNz[0], f_tilde, Ny)
res_exp0 = minimize(minimizer(Z0, f_tilde_Z0, 0, prob = 1), [1, 0.1, 0.01],
                bounds=((0.01, None), (0, None), (0, None)),
                method='L-BFGS-B')
res_SE0 = minimize(minimizer(Z0, f_tilde_Z0, 1, prob = 1), [1, 0.1, 0.001],
                bounds=((0, None), (0, None), (0, None)), tol = 1e-1,
                method='L-BFGS-B')

Z1, f_tilde_Z1, Y = create_Domain1(NNz[1], f_tilde, Ny)
res_exp1 = minimize(minimizer(Z1, f_tilde_Z1, 0, prob = 1), [1, 0.1, 0.01],
                bounds=((0.01, None), (0, None), (0, None)),
                method='L-BFGS-B')
res_SE1 = minimize(minimizer(Z1, f_tilde_Z1, 1, prob = 1), [1, 0.1, 0.001],
                bounds=((0, None), (0, None), (0, None)), tol = 1e-1,
                method='L-BFGS-B') # occhio ai parametri iniziali

# Compute a run with the optimal parameters
for KER in range(0,2):
    for PTS in range(0,2):

        #define the case study
        kernel, Z, f_tilde_Z, l_opt, var_opt, s_2_opt = case_study1(KER, PTS, Z0,
    f_tilde_Z0, Z1, f_tilde_Z1, res_exp0, res_SE0, res_exp1, res_SE1)

        #evaluate kernels
        K_Z, K_Y, K_ZY = evaluate_kernels(Z, Y, kernel, l_opt, var_opt)
```

```
57        #computes conditional mean and variance of the process
58        mu_YgivenZ, K_YgivenZ = prediction(f_tilde_Z, K_Z, K_Y, K_ZY, s_2_opt)
59        #samples again
60        y_cond=g_process(mu_YgivenZ, K_YgivenZ, n_process = 3)
```

## A.2 TASK2: GP REGRESSION ON A PERMABILITY FIELD

```
1  # load files and extract values
2  perm = np.load('true_perm.npy')
3
4  # define the coordinates of available observations
5  Coord1 = (np.array([np.linspace(15, 95, 5), np.linspace(10, 50, 5)]).T).astype(int)
6  Coord2 = (np.array([np.linspace(5, 105, 11), np.linspace(5, 55, 11)]).T).astype(int)
7
8  # define the domains
9  Mesh_total, Points_total = create_Domain2(perm)
10 Mesh1, Values1, Points1 = create_Dataset(perm, Coord1)
11 Mesh2, Values2, Points2 = create_Dataset(perm, Coord2)
12
13 # minimize the log-likelihood for Points1
14 res1_exp = minimize(minimizer(Points1, Values1, ker = 0, prob=2),
15              x0 = [1, 0.1, 0.1],
16              bounds=((1e-10, None), (1e-10, None), (1e-10, None)),
17              method='L-BFGS-B')
18 res1_se = minimize(minimizer(Points1, Values1, ker = 1, prob=2),
19               x0 = [5., 8, 1., 1.],
20               bounds=((1e-10, None), (1e-10, None), (1e-10, None),(1e-10, None)),
21               method='L-BFGS-B')
22
23 # minimize the log-likelihood for Points2
24 res2_exp = minimize(minimizer(Points2, Values2, ker = 0, prob=2),
25              x0 = [1, 0.1, 0.1],
26              bounds=((1e-10, None), (1e-10, None), (1e-10, None)),
27              method='L-BFGS-B')
28 res2_se = minimize(minimizer(Points2, Values2, ker = 1, prob=2),
29               x0 = [5., 10., 1., 1.],
30               bounds=((1e-10, None), (1e-10, None), (1e-10, None),(1e-10, None)),
31               method='L-BFGS-B')
32
33 # Compute a run with the optimal parameters
34
35 for KER in range(2):
36     for PTS in range(1,3):
37         kernel, Points, Values, l_opt, var_opt, s_2_opt = case_study2(KER, PTS,
    Points1, Values1, Points2, Values2, res1_exp, res1_se, res2_exp, res2_se)
38
39         # Fit Gaussian Process Regression
40         K_Z, K_Y, K_ZY = evaluate_kernels(Points, Points_total, kernel, l = l_opt,
    var = var_opt)
41         mean_perm, cov_perm = prediction(Values, K_Z, K_Y, K_ZY, s_2 = s_2_opt)
42         var_perm = np.diag(cov_perm)
43
44         # Generate gaussian processes with Circular Embedding
45         perm_cond = conditional(K_ZY, K_Z, K_Y, Points, Values, n_process = 1000)
46
47         # Evaluate the expected value of psi with Monte Carlo and Antithetic
    Variables
48         Critical = (np.array([np.ones(4)*85, np.array([35, 40, 45, 50])]).T).astype(
    int)
49         indexes = np.array([60*Critical[0,0] + i for i in Critical[:,1]])
```

```
50        mean_interest = mean_perm[indexes]
51        var_interest = cov_perm[indexes,:][:,indexes]
52        N, mean_est, std_est = MC_2stages(mean_interest, var_interest, tol = tol, N0
     = N0)
53
54        # Generate additional data for the convergence graph
55        Ns_all = np.array([100000, 250000, 500000, 1000000, 2500000, 5000000,
     10000000])
56        means_MC, stds_MC, means_AV, stds_AV = plot_MC_AV(Ns_all, mean_interest,
     var_interest, KER, PTS)
```

# B   APPENDIX: HELPERS

## B.1   KERNELS

```
1  # define euclidean-distance between m d-dim points
2  def dist_2(X1, X2):
3      rows, cols = np.indices((X1.shape[0], X2.shape[0]))
4      distances = np.sqrt(np.sum((X1[rows, :] - X2[cols, :])**2, axis=2))
5      return distances
6
7  # define weighted-distance between m d-dim points
8  def dist_w(X1, X2, l):
9      rows, cols = np.indices((X1.shape[0], X2.shape[0]))
10     distances = np.sum(((X1[rows, :] - X2[cols, :])/l)**2, axis=2)
11     return distances
12
13 # define exponential kernel
14 def kernel_exp(X1, X2, l, var):
15     distances = dist_2(X1, X2)
16     return var * np.exp(-distances/l)
17
18 # define square exponential kernel
19 def kernel_SE(X1, X2, l, var):
20     distances = dist_w(X1, X2, l)
21     return var * np.exp(-0.5*distances)
22
23 # evaluate kernels given the dataset and hyperparameters
24 def evaluate_kernels(Z, Y, kernel, l, var):
25     K_Z = kernel(Z,Z,l,var)
26     K_Y = kernel(Y,Y,l,var)
27     K_ZY = kernel(Z,Y,l,var)
28     return K_Z, K_Y, K_ZY
```

## B.2   DOMAINS

```
1  # create Dataset for task 1
2  def create_Domain1(Nz, f_tilde, Ny):
3      Z = 2*np.pi*np.random.uniform(size = Nz)
4      f_tilde_Z = f_tilde(Z)
5      Y = np.linspace(0, 2*np.pi, Ny).reshape(-1, 1)
6      return Z, f_tilde_Z, Y
7
8  # create datasets for task2
9  def create_Domain2(perm):
10     Y = perm.shape[0]
11     X = perm.shape[1]
12     Mesh_total = np.meshgrid(np.arange(Y), np.arange(X))
```

```
13      Points_total = np.array([Mesh_total[0].T.flatten(), Mesh_total[1].T.flatten()]).
    T
14      return Mesh_total, Points_total
15
16  def create_Dataset(perm, Coord):
17      Mesh = np.meshgrid(Coord[:,0], Coord[:,1])
18      Points = np.array([Mesh[0].T.flatten(), Mesh[1].T.flatten()]).T
19      Values = perm[Points[:, 0], Points[:,1]]
20      return Mesh, Values, Points
```

## B.3   OPTIMIZATION OF HYPERPARAMETERS

```
1   def minimizer(X, Y, ker, prob):
2       if ker == 0:
3           kernel = lambda x1, x2, l, var: kernel_exp(x1,x2,l,var)
4       elif ker == 1:
5           if prob == 1:
6               kernel = lambda x1, x2, l, var: kernel_SE(x1,x2,l,var)
7           elif prob == 2:
8               kernel = lambda x1, x2, l1, l2, var: kernel_SE(x1,x2,np.array([l1,l2]),
    var)
9       def log_likelihood(theta):
10          if (prob == 2 and ker == 1):
11              K = kernel(X, X, l1=theta[0], l2=theta[1], var=theta[2]) + theta[3] * np
    .eye(len(X))
12          else:
13              K = kernel(X, X, l=theta[0], var=theta[1]) + theta[2] * np.eye(len(X))
14
15          L = np.linalg.cholesky(K)
16          U = solve_triangular(L, Y, lower=True)
17          B = solve_triangular(L.T, U, lower=False)
18          logL = np.sum(np.log(np.diagonal(L))) + 0.5 * Y.dot(B)
19          return logL
20
21      return log_likelihood
```

## B.4   SAMPLING OF PROCESSES

```
1   # predict mean and covariance of the conditioned gaussian process
2   def prediction(f, K_Z, K_Y, K_ZY, s_2):
3       I = s_2*np.eye(K_Z.shape[0])
4       L = np.linalg.cholesky(K_Z+I)
5       mu_pred = K_ZY.T@(np.linalg.solve(L.T, np.linalg.solve(L, f)))
6       cov_pred = K_Y - K_ZY.T@(np.linalg.solve(L.T, np.linalg.solve(L, K_ZY)))
7       return mu_pred, cov_pred
8
9   # generate different gaussian processes
10  def g_process(mean, sigma, n_process = 1):
11      U,D,_ = np.linalg.svd(sigma)
12      A = U@np.diag(np.sqrt(D))
13      Ny = sigma.shape[0]
14      proc = mean + A@np.random.standard_normal(size = (Ny,n_process))
15      return proc
16
17  # Circulant Embedding in 2D
18  def build_P(sigma):
19      upper_left_P = np.empty((60,110))
20      for j in range(6600):
21          delta_y = int(j/60)
```

```python
22          delta_x = int(j%60)
23      upper_left_P[delta_x, delta_y] = sigma[0,j]
24      upper_right_P = np.flip(upper_left_P[:,1:], axis=1)
25      upper_P = np.c_[upper_left_P, upper_right_P]
26      lower_P = np.flip(upper_P[1:,:], axis=0)
27      P = np.r_[upper_P, lower_P]
28      return P
29
30  def CE(sigma, n_process):
31      P = build_P(sigma)
32      W = fft2(P)
33      gauss_process = np.zeros((6600, n_process))
34      for i in range(n_process):
35          Xi = np.random.standard_normal(size=W.shape)
36          T = 2*np.sqrt(59*109*W)*Xi
37          Z = ifft2(T)
38          Z_prime = Z.real + Z.imag
39          sample = (Z_prime[:60, :110]).T.flatten()
40          gauss_process[:,i] = sample
41      return gauss_process
42
43  # condition the zero-mean gaussian processes to the observations
44  def conditional(K_ZY, K_Z, K_Y, Points, Values, n_process):
45      index = 60*Points[:,0]+Points[:,1] #indexes of Z
46      precond = CE(K_Y, n_process)
47      z_iniz = precond[index]
48      post = precond + K_ZY.T@np.linalg.solve(K_Z, Values.reshape(-1,1)-z_iniz)
49      return post
```

### B.5 MONTE CARLO AND VARIANCE REDUCTION

```python
1  # define the function whose expected value we want to calculate
2  def csi(ip):
3      cnt = np.count_nonzero(ip>8, axis=0)
4      cnt[cnt>0]=1
5      return cnt
6
7  # define the Crude Monte Carlo method
8  def MC(mu, sigma, N):
9      ip = g_process(mu, sigma, n_process = N)
10     est = csi(ip)
11     mean_est = np.mean(est)
12     std_est = np.std(est, ddof = 1)
13     return mean_est, std_est
14
15 # define the 2stages Monte Carlo method
16 def MC_2stages(mu, sigma, tol, N0 = 100000):
17     alpha = 0.05
18     c = norm.ppf(1-alpha/2) #1.96
19     N_old = N0
20     mean_est_, std_est_ = MC(mu, sigma, N_old)
21     N = N_old
22     if(std_est_!=0):
23         N = int((c*std_est_/tol)**2)
24     mean_est, std_est = MC(mu, sigma, N)
25     while std_est >= std_est_:
26         N_old = N
27         mean_est_, std_est_ = MC(mu, sigma, N_old)
28         if std_est_ != 0:
29             N = int((c*std_est_/tol)**2)
```

```
30          mean_est, std_est = MC(mu, sigma, N)
31      return N, mean_est, std_est
32
33  # define the antithetic variables method
34  def MC_AV(mu, sigma, N):
35      N2 = int(N/2)
36      X = g_process(mu, sigma, n_process = N2)
37      Xa = 2*mu - X
38      Z = csi(X)
39      Za = csi(Xa)
40      est = (Z + Za)/2
41      mean_est = np.mean(est)
42      std_est = np.std(est, ddof=1)
43      return mean_est, std_est
44
45  # compute different simulations for MC and AV
46  def plot_MC_AV(Ns_all, mean_interest, var_interest, KER, PTS):
47      for i in range(len(Ns_all)):
48          means_MC[i], stds_MC[i] = MC(mean_interest, var_interest, Ns_all[i])
49          means_AV[i], stds_AV[i] = MC_AV(mean_interest, var_interest, Ns_all[i])
50      upperlim = means_MC + 1.96*stds_MC/np.sqrt(Ns_all)
51      lowerlim = means_MC - 1.96*stds_MC/np.sqrt(Ns_all)
52      upperlim = means_AV + 1.96*stds_AV/np.sqrt(Ns_all/2)
53      lowerlim = means_AV - 1.96*stds_AV/np.sqrt(Ns_all/2)
54      return means_MC, stds_MC, means_AV, stds_AV
```

# C  APPENDIX: GRADIENTS OF THE MARGINAL LOG-LIKELIHOOD

Before proceeding to calculate the marginal log-likelihood gradients analytically, we recall that, given a square matrix $A(\mathbf{x})$, the derivative of the matrix and the determinant result, respectively:

$$\frac{\partial A^{-1}}{\partial x_i} = -A^{-1}\frac{\partial A}{\partial x_i}A^{-1} \qquad \text{and} \qquad \frac{\partial \log|\det(A)|}{\partial x_i} = \text{Tr}\left(A^{-1}\frac{\partial A}{\partial x_i}\right)$$

Furthermore, for the sake of ease of notation we introduce $\tilde{\mathbf{K}} = \mathbf{K}(Z) + s^2\mathbf{I}$ and proceed with the analytical calculation:

- Notice that $\dfrac{\partial \tilde{\mathbf{K}}}{\partial s^2} = \mathbf{I}$, so:

$$\frac{\partial L}{\partial s^2} = -\frac{1}{2}v^T\left[-\tilde{\mathbf{K}}^{-1}\left(\frac{\partial \tilde{\mathbf{K}}}{\partial s^2}\right)\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\frac{\partial}{\partial s^2}\left[\log|\det(\tilde{\mathbf{K}})|\right] =$$

$$= \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\frac{1}{\det(\tilde{\mathbf{K}})}\frac{\partial \det(\tilde{\mathbf{K}})}{\partial s^2} =$$

$$= \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\text{Tr}\left(\tilde{\mathbf{K}}^{-1}\frac{\partial \tilde{\mathbf{K}}^{-1}}{\partial s^2}\right) =$$

$$= \frac{1}{2}v^T\left[\tilde{\mathbf{K}}^{-1}\tilde{\mathbf{K}}^{-1}\right]v - \frac{1}{2}\text{Tr}\left(\tilde{\mathbf{K}}^{-1}\right)$$

- one observes that $\dfrac{\partial \tilde{\mathbf{K}}}{\partial \sigma^2} = \dfrac{\partial \mathbf{K}}{\partial \sigma^2} = \mathbf{K}_{\sigma^2} = \begin{cases} \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x'}\|_2}{l}\right\} & (exp) \\ \exp\left\{-\frac{1}{2}\sum_{i=1}^{k}\frac{(x_i-x_i')^2}{\ell_i^2}\right\} & (se) \end{cases}$

$$\frac{\partial L}{\partial \sigma^2} = -\frac{1}{2}v^T \left[ -\tilde{\mathbf{K}}^{-1} \left( \frac{\partial \tilde{\mathbf{K}}}{\partial \sigma^2} \right) \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \frac{\partial}{\partial \sigma^2} \left[ \log |\det(\tilde{\mathbf{K}})| \right] =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\sigma^2} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \frac{1}{\det(\tilde{\mathbf{K}})} \frac{\partial \det(\tilde{\mathbf{K}})}{\partial \sigma^2} =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\sigma^2} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \operatorname{Tr} \left( \tilde{\mathbf{K}}^{-1} \frac{\partial \tilde{\mathbf{K}}^{-1}}{\partial \sigma^2} \right) =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\sigma^2} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \operatorname{Tr} \left( \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\sigma^2} \right)$$

- notice that $\dfrac{\partial \tilde{\mathbf{K}}}{\partial \ell_i} = \dfrac{\partial \mathbf{K}}{\partial \ell_i} = \mathbf{K}_{\ell_i} = \begin{cases} \frac{1}{\ell_i^2} \mathbf{K} & (exp) \\ \frac{1}{\ell_i^3} \mathbf{K} & (se) \end{cases}$

$$\frac{\partial L}{\partial \ell_i} = -\frac{1}{2}v^T \left[ -\tilde{\mathbf{K}}^{-1} \left( \frac{\partial \tilde{\mathbf{K}}}{\partial \ell_i} \right) \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \frac{\partial}{\partial \ell_i} \left[ \log |\det(\tilde{\mathbf{K}})| \right] =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\ell_i} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \frac{1}{\det(\tilde{\mathbf{K}})} \frac{\partial \det(\tilde{\mathbf{K}})}{\partial \ell_i} =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\ell_i} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \operatorname{Tr} \left( \tilde{\mathbf{K}}^{-1} \frac{\partial \tilde{\mathbf{K}}^{-1}}{\partial \ell_i} \right) =$$

$$= \frac{1}{2}v^T \left[ \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\ell_i} \tilde{\mathbf{K}}^{-1} \right] v - \frac{1}{2} \operatorname{Tr} \left( \tilde{\mathbf{K}}^{-1} \mathbf{K}_{\ell_i} \right)$$

## REFERENCES

[1] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*. Vol. 2. 3. MIT press Cambridge, MA, 2006.

[2] Fabio Nobile. *Stochastic Simulation: Lecture Notes*. 2022.

[3] Andreas Van Barel. 'Multilevel Monte Carlo Methods for Robust Optimization of Partial Differential Equations'. In: (2021).