

Bayesian linear models for large datasets: Markov chain Monte Carlo or Matheron's update rule

Hassan MAATOUK^(*)

(*) LAMPS, Université de Perpignan Via Domitia, 52 av. Paul Alduy,
66860 Cedex 9 Perpignan, France

hassan.maatouk@univ-perp.fr

JS2O 2-4 avril 2025



**Université
Perpignan**
Via Domitia

CRÉATRICE D'AVENIRS DEPUIS 1350



Outline I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Gaussian processes

The statistical problem of recovering an unknown function f from the training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with Gaussian noise is considered :

$$y_i = f(\mathbf{x}_i) + \underbrace{\epsilon_i}_{\text{noise}}, \quad \epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (1)$$

- f is an unknown function representing a physical phenomenon :

$$f : \mathcal{X} \subset \mathbb{R}^d \longrightarrow \mathbb{R}.$$

- The data $\mathbf{y} = f(\mathbf{X}) + \boldsymbol{\epsilon} \in \mathbb{R}^n$, with $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_n]^\top \in \mathbb{R}^n$.

Gaussian processes

The statistical problem of recovering an unknown function f from the training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with Gaussian noise is considered :

$$y_i = f(\mathbf{x}_i) + \underbrace{\epsilon_i}_{\text{noise}}, \quad \epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (1)$$

- ▶ f is an unknown function representing a physical phenomenon :

$$f : \mathcal{X} \subset \mathbb{R}^d \longrightarrow \mathbb{R}.$$

- ▶ The data $\mathbf{y} = f(\mathbb{X}) + \boldsymbol{\epsilon} \in \mathbb{R}^n$, with $\mathbb{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_n]^\top \in \mathbb{R}^n$.

Definition (Gaussian processes)

Let $(Y(\mathbf{x}))_{\mathbf{x} \in \mathcal{X}}$ be a stochastic process on \mathcal{X} such that $\mathbb{E}[Y(\mathbf{x})^2] < +\infty$. Then

- ▶ mean function : $\mu(\mathbf{x}) = \mathbb{E}[Y(\mathbf{x})]$;
- ▶ covariance function : $k(\mathbf{x}, \mathbf{x}') = \text{Cov}(Y(\mathbf{x}), Y(\mathbf{x}'))$.

A **Gaussian process (GP)** Y , i.e., $Y \sim \mathcal{GP}(\mu, k)$, is a stochastic process s.t.

$$[Y(\mathbf{x}_1), \dots, Y(\mathbf{x}_n)]^\top \sim \underbrace{\mathcal{N}(\underbrace{\boldsymbol{\mu}}_{\text{vector}}, \underbrace{\mathbf{K}}_{\text{matrix}})}_{\text{Gaussian vector}}.$$

Simulation of Gaussian processes

TABLE: Some popular covariance functions $k(x, x')$ used in Machine Learning.

| Name | Expression | Class |
|--------------------------|--|----------------------|
| Squared exponential (SE) | $\exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$ | \mathcal{C}^∞ |
| Matérn $\nu = 5/2$ | $\left(1 + \frac{\sqrt{5} x-x' }{\ell} + \frac{5(x-x')^2}{3\ell^2}\right) \exp\left(-\frac{\sqrt{5} x-x' }{\ell}\right)$ | \mathcal{C}^2 |
| Matérn $\nu = 3/2$ | $\left(1 + \frac{\sqrt{3} x-x' }{\ell}\right) \exp\left(-\frac{\sqrt{3} x-x' }{\ell}\right)$ | \mathcal{C}^1 |
| Exponential | $\exp\left(-\frac{ x-x' }{\ell}\right)$ | \mathcal{C}^0 |

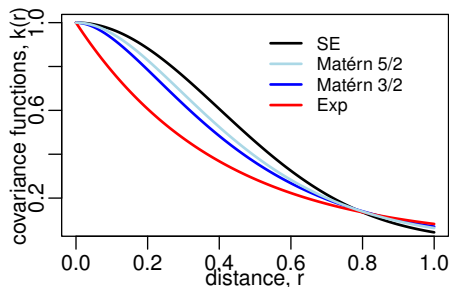


FIGURE: Some popular covariance functions (left) and GP sample paths using the Matérn $\nu = 3/2$ kernel (right).

Simulation of Gaussian processes

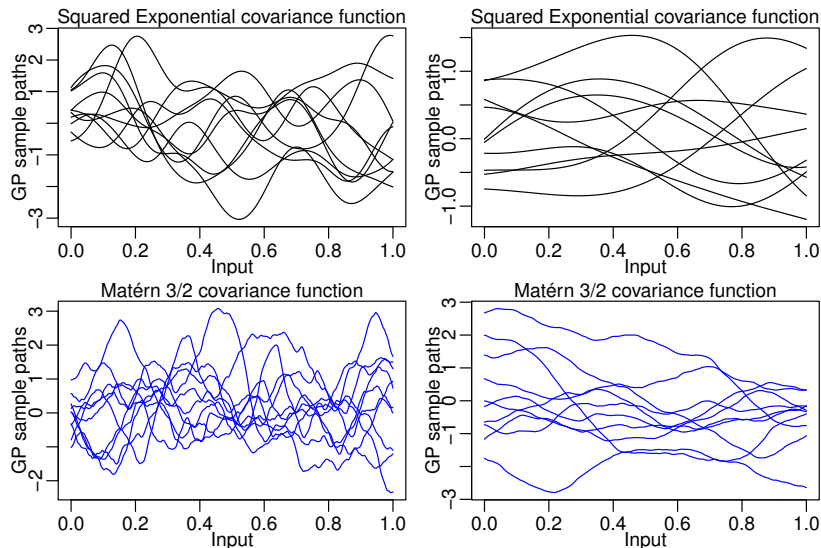


FIGURE: Ten zero-mean GP sample paths using the SE covariance function (top) and the Matérn covariance (bottom). The length-scale parameter ℓ is fixed at 0.1 (left) and at 0.4 (right).

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Gaussian process regression (GPR)

Gaussian process regression

- GPR is based on assuming a GP prior on the underlying function f .
- If $Y \sim \mathcal{GP}(0, k)$, then

$$\{Y|\mathbf{y}\} \sim \mathcal{GP}(\check{\mu}, \check{k}),$$

where the conditional mean $\check{\mu}$ and covariance function \check{k} are given as follows :

prediction $\rightarrow \quad \check{\mu}(x) = \mathbb{E}[Y(x)|\mathbf{y}] = k(x, \mathbb{X})^\top (k(\mathbb{X}, \mathbb{X}) + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{y};$

CI $\rightarrow \quad \check{k}(x, x') = k(x, x') - k(x, \mathbb{X})^\top (k(\mathbb{X}, \mathbb{X}) + \sigma^2 \mathbf{I}_n)^{-1} k(x', \mathbb{X}).$

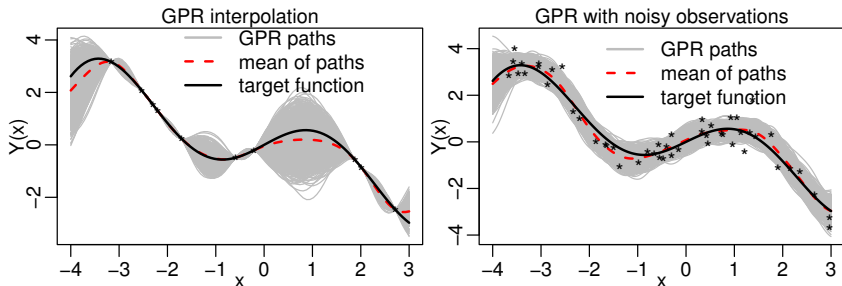


FIGURE: GPR : interpolation model ($\sigma = 0$, left) and noisy model ($\sigma = 0.5$, right).

Gaussian process regression (GPR)

Gaussian process regression

- GPR is based on assuming a GP prior on the underlying function f .
- If $Y \sim \mathcal{GP}(0, k)$, then

$$\{Y|\mathbf{y}\} \sim \mathcal{GP}(\check{\mu}, \check{k}),$$

where the conditional mean $\check{\mu}$ and covariance function \check{k} are given as follows :

prediction $\rightarrow \quad \check{\mu}(x) = \mathbb{E}[Y(x)|\mathbf{y}] = k(x, \mathbb{X})^\top (k(\mathbb{X}, \mathbb{X}) + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{y};$

CI $\rightarrow \quad \check{k}(x, x') = k(x, x') - k(x, \mathbb{X})^\top (k(\mathbb{X}, \mathbb{X}) + \sigma^2 \mathbf{I}_n)^{-1} k(x', \mathbb{X}).$

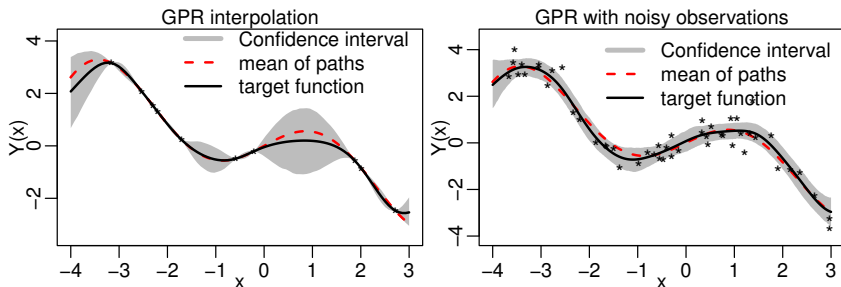


FIGURE: GPR : interpolation model ($\sigma = 0$, left) and noisy model ($\sigma = 0.5$, right).

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Bayesian linear models with large datasets

Definition (Finite-dimensional Bayesian linear models)

$(Y(\mathbf{x}))_{\mathbf{x} \in \mathcal{X}}$ is approximated by a finite-dimensional Bayesian linear model :

$$Y(\mathbf{x}) \approx \sum_{j=1}^N \xi_j \phi_j(\mathbf{x}) = \phi(\mathbf{x}) \boldsymbol{\xi} := Y^N(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X}, \quad (2)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{K})$ with $\tau^2 \mathbf{K}$ a positive-definite covariance matrix and $\phi(\cdot)$ is a deterministic function. Examples : KLE, Bernstein polynomials, B-splines, Hermite polynomials.

Bayesian linear models with large datasets

Definition (Finite-dimensional Bayesian linear models)

$(Y(x))_{x \in \mathcal{X}}$ is approximated by a finite-dimensional Bayesian linear model :

$$Y(x) \approx \sum_{j=1}^N \xi_j \phi_j(x) = \phi(x) \boldsymbol{\xi} := Y^N(x), \quad x \in \mathcal{X}, \quad (2)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{K})$ with $\tau^2 \mathbf{K}$ a positive-definite covariance matrix and $\phi(\cdot)$ is a deterministic function. Examples : KLE, Bernstein polynomials, B-splines, Hermite polynomials.

Bayesian linear models with large datasets

- Data : The dataset $\{Y^N(\mathbb{X}) + \epsilon = \mathbf{y}\}$ can be written in matrix form as follows :

$$\mathbf{X} \boldsymbol{\xi} + \epsilon = \mathbf{y},$$

where $\epsilon = [\epsilon_1, \dots, \epsilon_n]^\top$ is a zero-mean Gaussian noise vector with covariance matrix $\sigma^2 \mathbf{I}_n$, and $\mathbf{X} := \phi(\mathbb{X}) \in \mathbb{R}^{n \times N}$.

- Posterior distribution : Conditionally on the observations \mathbf{y}

$$\{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{K}}), \quad \text{where}, \quad (3)$$

$$\begin{cases} \tilde{\boldsymbol{\mu}} = \tau^2 (\mathbf{X} \mathbf{K})^\top (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{y}; \\ \tilde{\mathbf{K}} = \tau^2 \mathbf{K} - \tau^4 (\mathbf{X} \mathbf{K})^\top (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{X} \mathbf{K}. \end{cases} \quad (4)$$

Bayesian linear models with large datasets

Definition (Finite-dimensional Bayesian linear models)

$(Y(x))_{x \in \mathcal{X}}$ is approximated by a finite-dimensional Bayesian linear model :

$$Y(x) \approx \sum_{j=1}^N \xi_j \phi_j(x) = \phi(x) \boldsymbol{\xi} := Y^N(x), \quad x \in \mathcal{X}, \quad (2)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{K})$ with $\tau^2 \mathbf{K}$ a positive-definite covariance matrix and $\phi(\cdot)$ is a deterministic function. *Examples : KLE, Bernstein polynomials, B-splines, Hermite polynomials.*

Bayesian linear models with large datasets

- **Data** : The dataset $\{Y^N(\mathbb{X}) + \epsilon = \mathbf{y}\}$ can be written in matrix form as follows :

$$\mathbf{X} \boldsymbol{\xi} + \epsilon = \mathbf{y},$$

where $\epsilon = [\epsilon_1, \dots, \epsilon_n]^\top$ is a zero-mean Gaussian noise vector with covariance matrix $\sigma^2 \mathbf{I}_n$, and $\mathbf{X} := \phi(\mathbb{X}) \in \mathbb{R}^{n \times N}$.

- **Posterior distribution** : Conditionally on the observations \mathbf{y}

$$\{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{K}}), \quad \text{where}, \quad (3)$$

$$\begin{cases} \tilde{\boldsymbol{\mu}} = \tau^2 (\mathbf{X} \mathbf{K})^\top (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{y}; \\ \tilde{\mathbf{K}} = \tau^2 \mathbf{K} - \tau^4 (\mathbf{X} \mathbf{K})^\top (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{X} \mathbf{K}. \end{cases} \quad (4)$$

- The predictive equations in (4) require a matrix inversion of dimension $n \times n$, making this approach infeasible for a large number of observations n .

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Markov chain Monte Carlo (MCMC)

- According to [Williams and Rasmussen, 2006, Sect. 2.1.1] and Bayes' rule,

$$p(\boldsymbol{\xi}|\mathbf{X}, \mathbf{y}) := \frac{p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi})p(\boldsymbol{\xi})}{p(\mathbf{y}|\mathbf{X})}, \quad (5)$$

where $p(\mathbf{y}|\mathbf{X})$ is the normalizing constant, also known as the *marginal likelihood*. It is independent of $\boldsymbol{\xi}$ and given by

$$p(\mathbf{y}|\mathbf{X}) = \int_{\mathbb{R}^N} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi})p(\boldsymbol{\xi})d\boldsymbol{\xi}.$$

- By developing the likelihood $p(\mathbf{y}|\mathbf{X}, \mathbf{y})$ and the prior $p(\boldsymbol{\xi})$ in (5), we obtain

$$\begin{aligned} p(\boldsymbol{\xi}|\mathbf{X}, \mathbf{y}) &\propto \exp\left(-\frac{1}{2\sigma^2}[\mathbf{y} - \mathbf{X}\boldsymbol{\xi}]^\top[\mathbf{y} - \mathbf{X}\boldsymbol{\xi}]\right) \exp\left(-\frac{1}{2\tau^2}\boldsymbol{\xi}^\top \mathbf{K}^{-1}\boldsymbol{\xi}\right) \\ &\propto \exp\left(-\frac{1}{2}[\boldsymbol{\xi} - \boldsymbol{\mu}]^\top \boldsymbol{\Sigma}^{-1}[\boldsymbol{\xi} - \boldsymbol{\mu}]\right), \end{aligned} \quad (6)$$

where

$$\begin{cases} \boldsymbol{\mu} = \left[\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2\right]^{-1} \mathbf{X}^\top \mathbf{y}/\sigma^2; \\ \boldsymbol{\Sigma} = \left[\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2\right]^{-1}. \end{cases} \quad (7)$$

Markov chain Monte Carlo (MCMC)

Markov chain Monte Carlo

The posterior pdf in (6) is proportional to the product of a likelihood function and a zero-mean Gaussian :

$$\begin{aligned} p(\boldsymbol{\xi}|\mathbf{X}, \mathbf{y}) &\propto \underbrace{\exp\left(-\frac{1}{2\sigma^2}[\mathbf{y} - \mathbf{X}\boldsymbol{\xi}]^\top[\mathbf{y} - \mathbf{X}\boldsymbol{\xi}]\right)}_{L(\boldsymbol{\xi}): \text{likelihood function}} \underbrace{\exp\left(-\frac{1}{2\tau^2}\boldsymbol{\xi}^\top \mathbf{K}^{-1}\boldsymbol{\xi}\right)}_{\text{Gaussian prior}} \\ &\propto L(\boldsymbol{\xi})\mathcal{N}(\boldsymbol{\xi}; \mathbf{0}, \tau^2 \mathbf{K}). \end{aligned} \quad (8)$$

The logarithm function in (8), which has a computational complexity of order $\mathcal{O}(nN)$ will be evaluated at each MCMC iteration.

- Sampling : Metropolis-Hastings (MH) proposals [Neal, 1999] :

$$\boldsymbol{\xi}' = \rho \boldsymbol{\nu} + \sqrt{1 - \rho^2} \boldsymbol{\xi}, \quad \boldsymbol{\nu} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{K}), \quad (9)$$

where $\rho \in [-1, 1]$ is a step-size parameter, $\boldsymbol{\xi}$ is the current state, and $\boldsymbol{\xi}'$ is the proposal state.

Recall that the MH acceptance ratio, $\alpha = \min \{1, L(\boldsymbol{\xi}')/L(\boldsymbol{\xi})\}$ depends solely on the likelihood ratio and is independent of ρ .

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Matheron's update rule

Proposition (Matheron's update rule (MUR))

Let $\boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \tau^2 \mathbf{K})$. Suppose that $\mathbf{X} \in \mathbb{R}^{n \times N}$ is a given matrix of rank n , and $\mathbf{y} \in \mathbb{R}^n$ is an output vector representing the data i.e., $\mathbf{X}\boldsymbol{\xi} = \mathbf{y}$. Then

$$\{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \underbrace{\boldsymbol{\xi}}_{\text{prior}} + \underbrace{(\mathbf{X}\mathbf{K})^\top (\mathbf{X}\mathbf{K}\mathbf{X}^\top)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\xi})}_{\text{update}}. \quad (10)$$

Additionally, we have

$$\phi(\cdot) \{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \phi(\cdot) \left[\boldsymbol{\xi} + (\mathbf{X}\mathbf{K})^\top (\mathbf{X}\mathbf{K}\mathbf{X}^\top)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\xi}) \right],$$

where $\phi(\cdot)$ is the basis vector appearing in the Bayesian linear model.

Matheron's update rule

Proposition (Matheron's update rule (MUR))

Let $\xi \sim \mathcal{N}(\mu, \tau^2 \mathbf{K})$. Suppose that $\mathbf{X} \in \mathbb{R}^{n \times N}$ is a given matrix of rank n , and $\mathbf{y} \in \mathbb{R}^n$ is an output vector representing the data i.e., $\mathbf{X}\xi = \mathbf{y}$. Then

$$\{\xi | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \underbrace{\xi}_{\text{prior}} + \underbrace{(\mathbf{X}\mathbf{K})^\top (\mathbf{X}\mathbf{K}\mathbf{X}^\top)^{-1} (\mathbf{y} - \mathbf{X}\xi)}_{\text{update}}. \quad (10)$$

Additionally, we have

$$\phi(\cdot) \{\xi | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \phi(\cdot) \left[\xi + (\mathbf{X}\mathbf{K})^\top (\mathbf{X}\mathbf{K}\mathbf{X}^\top)^{-1} (\mathbf{y} - \mathbf{X}\xi) \right],$$

where $\phi(\cdot)$ is the basis vector appearing in the Bayesian linear model.

MUR with Noisy observations

More generally, if the data are observed with independent Gaussian noise $\{\mathbf{X}\xi + \epsilon = \mathbf{y}\}$, where $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, then

$$\{\xi | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \xi + \underbrace{\tau^2 (\mathbf{X}\mathbf{K})^\top (\tau^2 \mathbf{X}\mathbf{K}\mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1}}_{n \times n \text{ matrix}} (\mathbf{y} - \mathbf{X}\xi - \epsilon). \quad (11)$$

Matheron's update rule

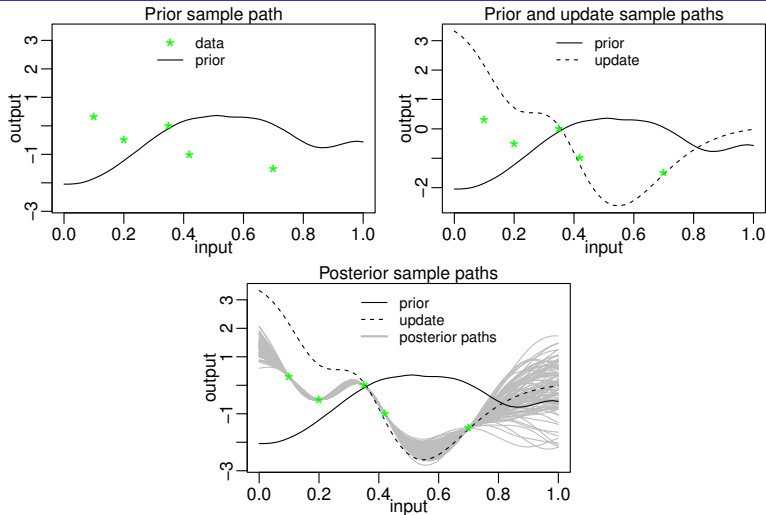


FIGURE: Visual representation of the MUR (*noise-free case*). **Top left** : a single path of the prior together with the data (green stars). **Top right** : the corresponding update sample path, derived from (10) is displayed as black dashed curve. **Bottom** : posterior paths (gray solid curves) are obtained by combining the prior and the update as per (10).

Matheron's update rule for large datasets

Proposition (MUR for large datasets [Maatouk et al., 2025])

Under the same settings as Proposition 1, we have

$$\{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} \stackrel{d}{=} \underbrace{\boldsymbol{\xi}}_{\text{prior}} + \underbrace{(\mathbf{X}^\top \mathbf{X} / \sigma^2 + \mathbf{K}^{-1} / \tau^2)^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\xi} - \boldsymbol{\epsilon}) / \sigma^2}_{\text{update}},$$

where $\mathbf{y} - \mathbf{X} \boldsymbol{\xi} - \boldsymbol{\epsilon}$ represents the residual and σ^2 is the variance of the noise.

Comments

- The update part requires a matrix inversion of dimension $N \times N$ instead of $n \times n$.
- As for the MCMC approaches, sampling is performed before conditioning rather than after.
- Unlike MCMC approaches, we do not need to evaluate a likelihood function at each iteration.

Matheron's update rule for large datasets

Before proving Proposition 2, let us present the following Lemma.

Lemma

Consider three random vectors $\mathbf{V}_1 \in \mathbb{R}^N$, $\mathbf{V}_2 \in \mathbb{R}^n$ and $\mathbf{V}_3 \in \mathbb{R}^N$ s.t.

$$\mathbf{V}_1 \stackrel{d}{=} f(\mathbf{V}_2) + \mathbf{V}_3,$$

where f is a measurable function of \mathbf{V}_2 and where \mathbf{V}_2 is independent of \mathbf{V}_3 .
Then,

$$\{\mathbf{V}_1 | \mathbf{V}_2 = \boldsymbol{\theta}\} \stackrel{d}{=} f(\boldsymbol{\theta}) + \mathbf{V}_3,$$

for any $\boldsymbol{\theta} \in \mathbb{R}^n$.

Démonstration.

The proof is provided in [Wilson et al., 2021, Lemma 2].



Matheron's update rule for large datasets

Proof of Proposition 2 Maatouk, Rulli re and Bay (2025).

From the equivalent between the two *direct* approaches Equations (4) and (7), we have

$$(\mathbf{X}^\top \mathbf{X} / \sigma^2 + \mathbf{K}^{-1} / \tau^2) \mathbf{X}^\top / \sigma^2 = \tau^2 \mathbf{K} \mathbf{X}^\top (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1}. \quad (12)$$

Let $\mathbf{V}_3 := \boldsymbol{\xi} - (\mathbf{X}^\top \mathbf{X} / \sigma^2 + \mathbf{K}^{-1} / \tau^2)^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}) / \sigma^2$. Additionally, we have

$$\mathbb{E}[\boldsymbol{\xi} | \mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}] = (\mathbf{X}^\top \mathbf{X} / \sigma^2 + \mathbf{K}^{-1} / \tau^2)^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}) / \sigma^2.$$

Thus, we can write :

$$\begin{aligned} \boldsymbol{\xi} &= \mathbb{E}[\boldsymbol{\xi} | \mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}] + (\boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi} | \mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}]) \\ &= (\mathbf{X}^\top \mathbf{X} / \sigma^2 + \mathbf{K}^{-1} / \tau^2)^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}) / \sigma^2 + \mathbf{V}_3. \end{aligned}$$

Let $\mathbf{V}_1 = \boldsymbol{\xi}$ and $\mathbf{V}_2 = \mathbf{X} \boldsymbol{\xi} + \boldsymbol{\epsilon}$. Since \mathbf{V}_2 and \mathbf{V}_3 are jointly Gaussian but uncorrelated, it follows that they are independent. □

Matheron's update rule for large datasets

Rest of the proof.

Indeed,

$$\begin{aligned}\text{Cov}(\mathbf{V}_2, \mathbf{V}_3) &= \text{Cov}(\mathbf{X}\boldsymbol{\xi} + \boldsymbol{\epsilon}, \boldsymbol{\xi} - (\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\xi} + \boldsymbol{\epsilon})/\sigma^2) \\ &= \tau^2 \mathbf{X} \mathbf{K} - \text{Var}(\mathbf{X}\boldsymbol{\xi} + \boldsymbol{\epsilon})[(\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top / \sigma^2]^\top \\ &= \tau^2 \mathbf{X} \mathbf{K} - (\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)(\tau^2 \mathbf{X} \mathbf{K} \mathbf{X}^\top + \sigma^2 \mathbf{I}_n)^{-1} \tau^2 \mathbf{X} \mathbf{K} \\ &= \tau^2 \mathbf{X} \mathbf{K} - \tau^2 \mathbf{X} \mathbf{K} = \mathbf{0}_{n,N},\end{aligned}$$

where $\mathbf{0}_{n,N}$ is the $n \times N$ zero matrix. The second-to-last line is done using Equation (12). Setting $f(\mathbf{V}_2) = (\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top \mathbf{V}_2/\sigma^2$ and using Lemma 1, we obtain

$$\begin{aligned}\{\boldsymbol{\xi} | \mathbf{X}, \mathbf{y}\} &\stackrel{d}{=} (\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top \mathbf{y}/\sigma^2 + \boldsymbol{\xi} - \\ &\quad (\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\xi} + \boldsymbol{\epsilon})/\sigma^2 \\ &\stackrel{d}{=} \boldsymbol{\xi} + (\mathbf{X}^\top \mathbf{X}/\sigma^2 + \mathbf{K}^{-1}/\tau^2)^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\xi} - \boldsymbol{\epsilon})/\sigma^2.\end{aligned}$$

Hence, the claim follows. □

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
 - Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Illustrative examples 1D

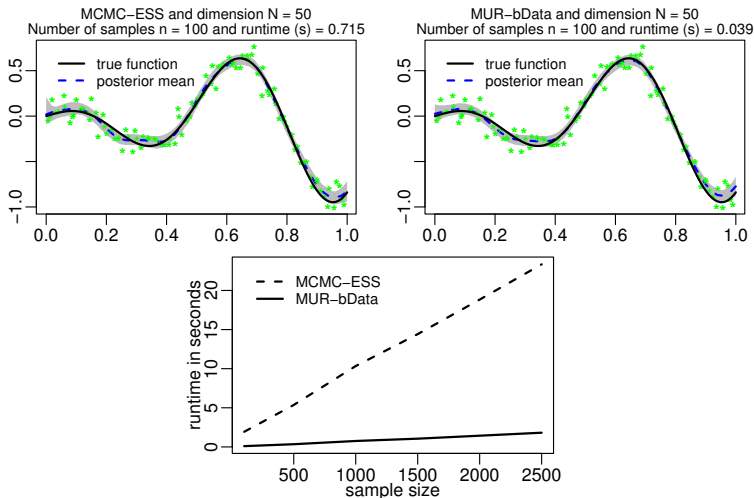


FIGURE: Top : Performance accuracy of the finite-dimensional Bayesian linear for $N = 50$. The number of samples (green stars) is fixed at $n = 100$. The gray shaded area represents the 95% confidence interval based on **6,000** sample paths. The highly efficient MCMC approach is employed in the left panel, while the proposed MUR for big data is applied in the right panel. **Bottom** : Runtime in seconds for generating **15,000** posterior sample paths as a function of the number of samples for the two competing approaches.

Illustrative examples 1D (big data and extreme case)

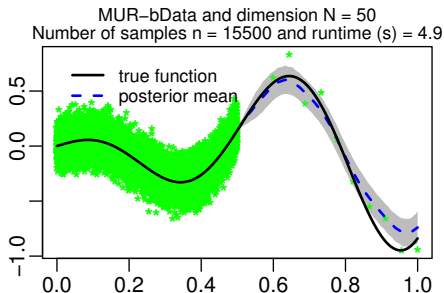
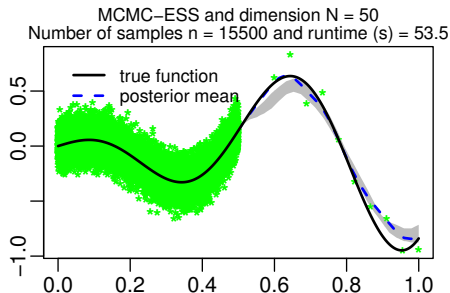


FIGURE: Same settings as Fig. 5 except the number of samples n which is fixed at **15,500** instead of **100**. There are **15,490** observations available in the first half of the domain, while only **10** are available in the second. Unlike MUR-bData, the 95% confidence interval (gray shaded area) of the MCMC-ESS approach does not closely follow the posterior mean (blue dashed curve).

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
 - Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Real-world diamond data

Now, we apply the two strategies developed in this presentation to real-world diamond data. This dataset consists of the prices in US dollars (326\$-18,823\$) of $n = 53,940$ diamonds as a function of their carat, i.e., weight of the diamond (0.2-5.01).

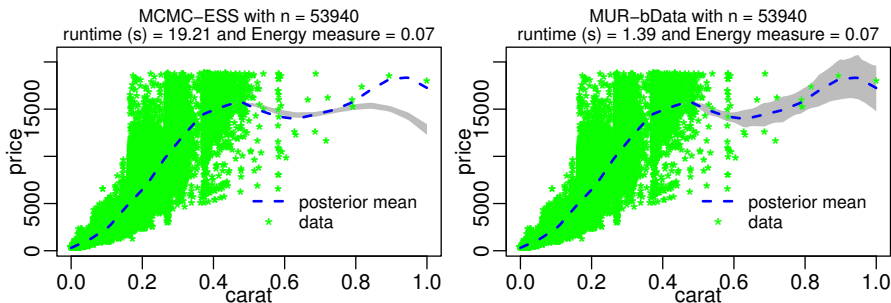


FIGURE: Accuracy estimation of the price of $n = 53,940$ diamonds as a function of carat. The two developed strategies are employed MCMC-ESS (left panel) and MUR-bData (right panel). The computational running time of generating 1,000 sample paths and the energy measure criterion are displayed in the main of each panel.

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Illustrative examples 2D

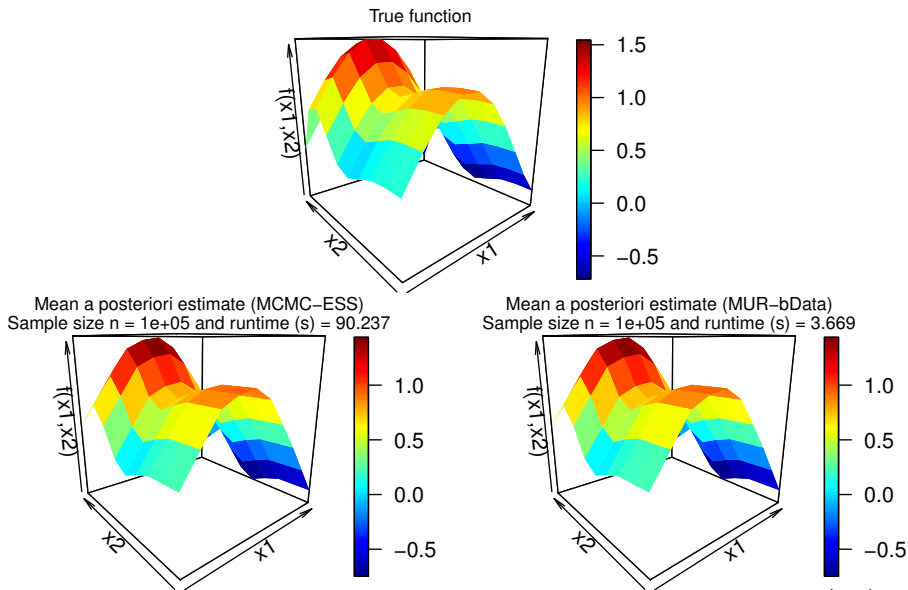


FIGURE: Computational running time for generating 2,000 surfaces using MCMC (left) and MUR-bData (right) on large datasets. The number of observations is $n = 100,000$.

Table of Contents I

1 Gaussian process regression review

- Gaussian processes
- Gaussian process regression

2 Bayesian linear models with large datasets

- Introduction
- First approach : Markov chain Monte Carlo (MCMC)
- Second approach : Matheron's update rule (MUR)

3 Performance illustrations

- Illustrative examples 1D
- Real-world application : Diamond data
- Illustrative examples 2D

4 Conclusion

Conclusion

| Characteristic | MCMC-ESS | MUR-bData |
|----------------|------------------------------|--------------------------------|
| | Approximate approach | Exact approach |
| | Sampling before conditioning | Sampling before conditioning |
| | Fast | Faster |
| | Flexible approach | Less flexible approach |
| | Iterative simulation | Iterative or direct simulation |

TABLE: Comparison between MCMC-ESS and MUR-bData approaches.

References I



Maatouk, H., Rulli re, D., and Bay, X. (2025).

Bayesian linear models for large datasets : Markov chain Monte Carlo or Matheron's update rule.

Under review in MCQMC2024.



Neal, R. (1999).

Regression and classification using Gaussian process priors.

Bayesian Statistics, 6 :475–501.



Williams, C. K. and Rasmussen, C. E. (2006).

Gaussian processes for machine learning, volume 2.

MIT press Cambridge, MA.



Wilson, J., Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth, M. (2021).

Pathwise conditioning of Gaussian processes.

JMLR, 22(105) :1–47.

QUESTIONS ?

Thank you for your attention.