

On Multi-Output Kriging and Constrained Classification

Didier Rullière, Mines Saint-Etienne, LIMOS, drulliere@emse.fr
joint work with Marc Grossouvre, Mines Saint-Etienne, LIMOS, URBS

Séminaire de Statistique, Laboratoire de Mathématiques d'Avignon (LMA),
October 02, 2023



picture: mining headframe (chevalement) at Saint-Etienne

preliminary note

preliminary note

Very recent work (<3 months)

all comments are welcome 😊

Outline

- 1 Kriging: a brief overview
- 2 Constrained Multi-Output Kriging
- 3 Constrained Classification
- 4 Numerical Illustrations

Kriging: a brief overview

The Origins of Kriging

Geostatistical problem

- How to predict **gold concentration** everywhere in the ground,
- ... if it is only measured on few specific sites?

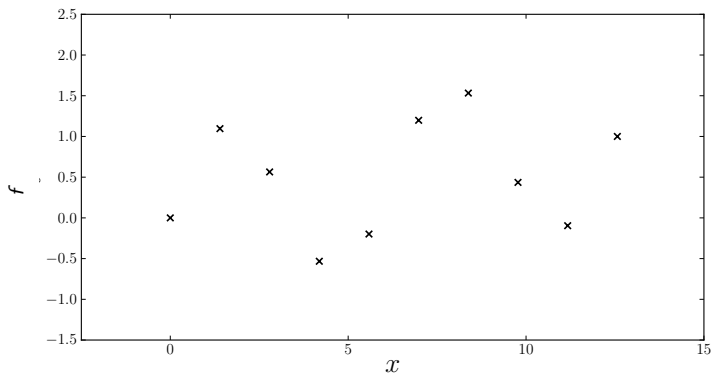


- This is an **interpolation** problem.

Gaussian Process Regression (1/5)

Gaussian approach: Gaussian Process Regression (\pm ML community)

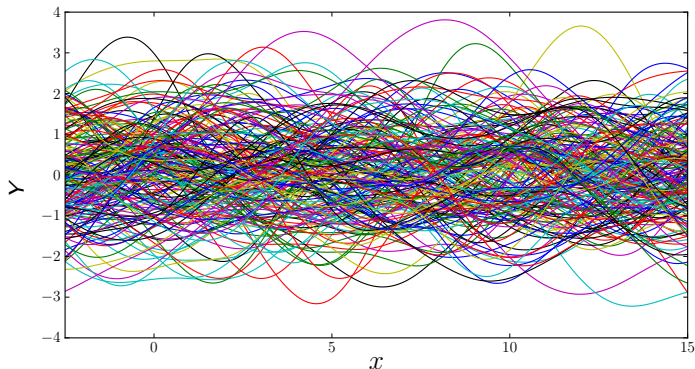
Assume we have observed a function $f(\cdot)$ over a set of points $\mathbb{X} = (x_1, \dots, x_n)^\top$:



Here x in 1D, f in 1D. The vector of observations is $\mathbf{y} = f(\mathbb{X})$, i.e. $y_i = f(x_i)$.

Gaussian Process Regression (2/5)

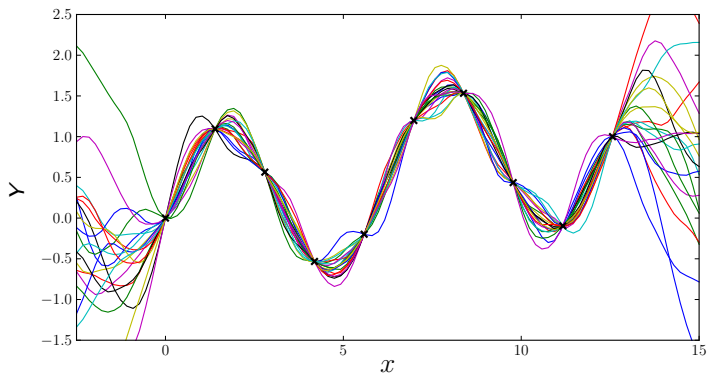
Since $f(\cdot)$ is unknown, we make the general assumption that it is close to the sample path of a Gaussian process $Y \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$:



here $\mu(x) = 0$.

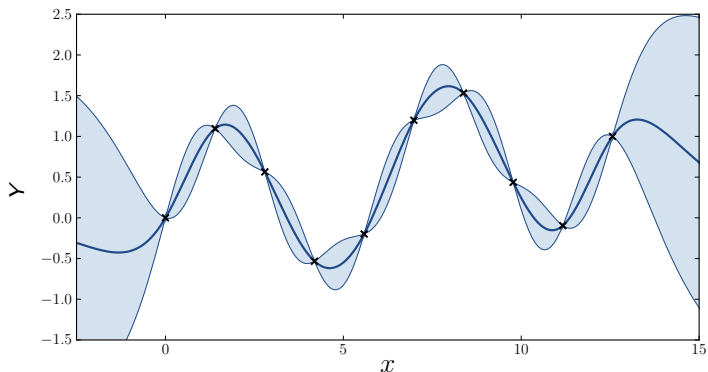
Gaussian Process Regression (3/5)

If we remove all the samples that do not interpolate the observations we obtain:



Gaussian Process Regression (4/5)

It can be summarized by a mean function and 95% confidence intervals.



- **Kriging Mean:** blue thick line
- **Kriging Standard Deviation:** proportional to confidence band width.

You can play here: https://durrande.shinyapps.io/gp_playground/ thanks Nicolas!

... here $x \in \mathbb{R}$, but it is easy to extend to $x \in \mathbb{R}^d$...

Gaussian Process Regression (5/5): equations

The conditional distribution of $Y(x^*)$ given $\mathbf{Y} = (Y(x_1), \dots, Y(x_n))^T$ can be obtained:

By definition, $(Y(x^*), \mathbf{Y})$ is multivariate normal. Formulas on the conditioning of Gaussian vectors give the distribution of $Y(x^*)|\mathbf{Y} = \mathbf{y}$. It is $\mathcal{N}(m(\cdot), c(\cdot, \cdot))$ with :

$$\begin{cases} m(x^*) &= \mu(x^*) + k(x^*, \mathbb{X})k(\mathbb{X}, \mathbb{X})^{-1}(\mathbf{y} - \mu(\mathbb{X})) \\ c(x^*, x^{*'}) &= k(x^*, x^{*'}) - k(x^*, \mathbb{X})k(\mathbb{X}, \mathbb{X})^{-1}k(\mathbb{X}, x^{*'}) \end{cases}$$

where $k(\cdot, \cdot)$ is the covariance function of the prior Gaussian Process.

Simple Kriging, Gaussian case

For a centered process, when $\mu(x^*) = 0$ for all x^* , the simple Kriging predictor mean and variance are

$$\begin{cases} \mathbb{E}[Y(x^*)|\mathbf{Y}=\mathbf{y}] &= m(x^*) &= \mathbf{h}(x^*)^T \mathbb{K}^{-1} \mathbf{y} \\ \text{Var}[Y(x^*)|\mathbf{Y}=\mathbf{y}] &= c(x^*, x^*) &= \sigma(x^*)^2 - \mathbf{h}(x^*)^T \mathbb{K}^{-1} \mathbf{h}(x^*) \end{cases}$$

with $\mathbb{K} = k(\mathbb{X}, \mathbb{X})$, $\mathbf{h}(x^*)^T = k(x^*, \mathbb{X})$ and $\sigma(x^*)^2 = k(x^*, x^*)$

Simple Kriging: the statistical approach (1/2)

Another Approach: Best Linear Unbiased Prediction (\pm Geostat community)

Define the linear predictor

$$M(x^*) := \sum_{i=1}^n \alpha_i(x^*) Y(x_i) = \boldsymbol{\alpha}(x^*)^\top \mathbf{Y}.$$

Now let us minimize on $\boldsymbol{\alpha}(x^*) = (\alpha_1(x^*), \dots, \alpha_n(x^*))$ the loss

$$\begin{aligned} \Delta(x^*) &:= \mathbb{E} \left[(M(x^*) - Y(x^*))^2 \right] \\ &= \boldsymbol{\alpha}(x^*)^\top \mathbb{K} \boldsymbol{\alpha}(x^*) - 2\boldsymbol{\alpha}(x^*)^\top \mathbf{h}(x^*) + \text{constant}. \end{aligned}$$

This leads to the vector of weights

$$\boldsymbol{\alpha}(x^*) = \mathbb{K}^{-1} \mathbf{h}(x^*),$$

where $\mathbf{h}(x^*)$ is the covariance vector between $Y(x^*)$ and the vector \mathbf{Y} , and \mathbb{K} is the covariance matrix of \mathbf{Y} . $Y(\cdot)$ centered here.

Simple Kriging: the statistical approach (2/2)

Predictor and variance

From that follows the expression of $M(x^*)$ and $\Delta(x^*)$:

$$\begin{cases} M(x^*) &= \mathbf{h}(x^*)^\top \mathbb{K}^{-1} \mathbf{Y} \\ \Delta(x^*) &= \sigma(x^*)^2 - \mathbf{h}(x^*)^\top \mathbb{K}^{-1} \mathbf{h}(x^*) \end{cases}$$

One retrieves exactly Kriging mean and variance. 😊

Notice that $\Delta(x^*)$ does not depend on observed responses \mathbf{Y} .

Pros and Cons of both approaches

- GPR more intuitive for varying x^*
- GPR more suited to Bayesian analysis and interpretation, more visual
- Stat Approach not limited to Gaussian case
- Stat Approach easier to extend (other combination, criterions, penalization, etc.)

Because GPR/Kriging predicts a full, spatially varying, distribution,
it is of great use in decision making.


⇒ a stepping stone to address decision-theoretic problems.

GPR/Kriging Problems

Here are few selected problems:

- **Model selection:**
how to choose prior process, prior covariance function family, prior covariance function parameters ?
- **Computation:**
how to compute the predictor when the matrices are huge?
- **Adaptation:**
how to adapt to specific settings (monotony, uncertainty, extremes, high dimension, **multiple outputs, constraints...**)

We focus here on the last problems: **multiple outputs, constraints**.

Main differences with usual Kriging models are highlighted with a symbol  .

Constrained Multi-Output Kriging

Constrained Multi-Output Kriging: motivation

Why another study, why multiple outputs, why specific constraints?

- **Multi-output:**

Studying multiple outputs is useful:

- Observations of $p > 1$ variables, possibly dependent
- Need for a model with not too many hyperparameters, not $\mathcal{O}(p^2)$

- **Constraints:**

Prescribing e.g. the average value of predictions is useful:

- external information (known quantity of chemical loss, national statistic...)
- adverse modelling (regulation, simulation under specific scenarios...)
- need to homogenize results (over different regions, observed years, fairness constraints...)

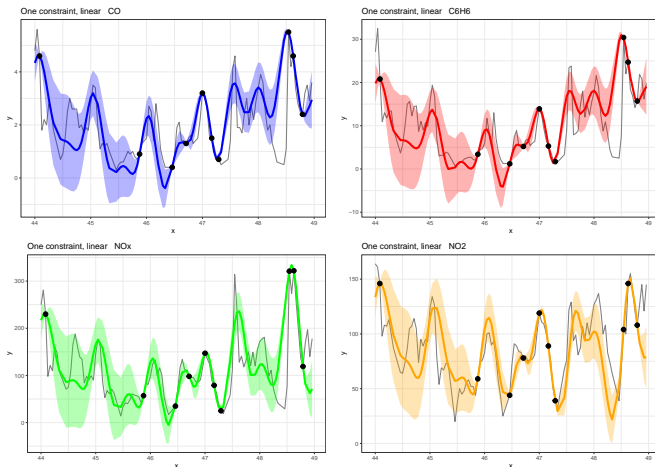
- **Classification in mind:**

Easy to adapt to classification, and useful:

- Multi-output applied to membership degrees
- Useful constraints: membership degrees sum to 1, prescribed percentages of each class.

Multi-Output Kriging

Example of multi-output, e.g. at $x \simeq 44$, an observation $\mathbf{Y}(x) = (Y_1(x), \dots, Y_4(x))^T$



Here time serie, $x \in \mathbb{R}$ for visual illustration, but in general $x \in \mathbb{R}^d$.

Literature

Among a vast literature,

- **Co-Kriging techniques:** usually one main “primary” output, $\mathcal{O}(p^2)$ covariance models, cross co-variograms *Goovaerts, 1998, Ver Hoef and Cressie, 1993, Furrer and Genton, 2011.*
- **Indicator Kriging:** with a latent process, post-treatments, many covariances models, *Journel, 1983, Meer, 1996, Goovaerts, 2009, Chiang et al., 2013, Agarwal et al., 2021*
- **Gaussian Process and classification:** with latent GP, Bayesian inference and approximations, *Williams and Barber, 1998, Rasmussen et al., 2006, Dahl and Bonilla, 2019, Panos et al., 2021.*
- **Constraints:** without Kriging, classification *Gordon, 1996, Bradley et al., 2000, Höppner and Klawonn, 2008, Ganganath et al., 2014*, fuzzy classification *Benatti et al., 2022*, fairness constraints *Zafar et al., 2019*

Multi-Output Kriging

Framework

- **Inputs:** set of locations $x \in \mathcal{X} \subset \mathbb{R}^d$.
- **Outputs:** multi-valued random field $\mathbf{Y}(x) := (Y_1(x), \dots, Y_p(x))^T \in \mathbb{R}^p$.
- **Observations:** $\mathbf{Y}(x_1), \dots, \mathbf{Y}(x_n)$.

Question

How to predict $\mathbf{Y}(\cdot)$ at some unobserved locations x_1^*, \dots, x_q^* ?

Joint Kriging Model

$$\mathbf{M}(x^*) := \sum_{i=1}^n \alpha_i(x^*) \mathbf{Y}(x_i), \quad \alpha_i(x^*) \in \mathbb{R}, \quad i = 1 \dots n \quad (1)$$

Main assumption: the weights are impacting all components the same way. 

Constrained Multi-Output Kriging

Recall the joint Kriging model,

$$\mathbf{M}(x^*) := \sum_{i=1}^n \alpha_i(x^*) \mathbf{Y}(x_i) = \mathbf{Y} \boldsymbol{\alpha}(x^*) \quad (2)$$

where $\mathbf{Y} := [\mathbf{Y}(x_1), \dots, \mathbf{Y}(x_p)] \in \mathbb{R}^{p \times n}$ and $\boldsymbol{\alpha}(x^*) := (\alpha_1(x^*), \dots, \alpha_n(x^*))^\top \in \mathbb{R}^n$.

How to get optimal weights?

To get optimal weights $\mathbb{A} := [\alpha(x_1^*), \dots, \alpha(x_q^*)]$, they are optimized in order to:

- minimize some error:

$$\Delta(x^*) := \mathbb{E} [\|\mathbf{M}(x^*) - \mathbf{Y}(x^*)\|_{\mathbb{W}}^2] \in \mathbb{R}. \quad \text{NEW} \quad (3)$$

where $\|\mathbf{v}\|_{\mathbb{W}}^2 := \mathbf{v}^\top \mathbb{W} \mathbf{v}$ and \mathbb{W} a given real symmetrical positive-definite matrix.

- under various constraints:

- **Constraint 1:** Sum of weights equal to 1, $\alpha_1(x^*) + \dots + \alpha_n(x^*) = 1$

- **Constraint 2:** Prescribed average \mathbf{m} of predicted values $\mathbf{M}(x_1^*), \dots, \mathbf{M}(x_q^*)$ NEW

Optimal weights, no constraint

Without constraint one retrieves Simple Kriging equations, but \mathbb{K} , $\mathbf{h}(\cdot)$ involve cross covariances of all components of $\mathbf{Y}(\cdot)$.

Proposition (Simple Joint Kriging weights)

The optimal weights $\alpha(x^*)$ minimizing the loss $\Delta(x^*)$ are given by:

$$\alpha(x^*) = \mathbb{K}^{-1}\mathbf{h}(x^*), \quad (4)$$

or equivalently, using a matrix expression, under invertibility assumption,

$$\mathbb{A} = \mathbb{K}^{-1}\mathbb{H}, \quad (5)$$

where $\mathbb{K} := \mathbb{E}[\mathbf{Y}^\top \mathbb{W} \mathbf{Y}]$, $\mathbf{h}(x^*) := \mathbb{E}[\mathbf{Y}^\top \mathbb{W} \mathbf{Y}(x^*)]$, and $\mathbb{H} := [\mathbf{h}(x_1^*), \dots, \mathbf{h}(x_q^*)]$. 

If furthermore $\mathbb{E}[Y_j(x)] = 0$ for all $j = 1, \dots, p$, $x \in \mathcal{X}$, then $\mathbf{M}(x^*)$ is unbiased.

matrix sizes: $\alpha(x^*) \in \mathbb{R}$, $\mathbb{K} \in \mathbb{R}^{n \times n}$, $\mathbf{h}(x^*) \in \mathbb{R}^n$, $\mathbb{H} \in \mathbb{R}^{n \times q}$.

Optimal weights, constraint 1

Considered constraint, similarly to ordinary Kriging

Constraint 1: Weights sum to one, $\alpha^\top(x^*)\mathbf{1}_n = 1$, $x^* \in \mathcal{X}$.

Proposition (Ordinary Joint Kriging weights)


Under the Constraint 1, the optimal weights $\alpha(x^)$ minimizing the loss $\Delta(x^*)$ are:*

$$\begin{cases} \alpha(x^*) &= \mathbb{K}^{-1}(\mathbf{h}(x^*) + \lambda(x^*)\mathbf{1}_n) \\ \lambda(x^*) &= \frac{1}{\delta}(\mathbf{1} - \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbf{h}(x^*)) \end{cases} \quad (6)$$

Equivalently, using matrix expressions, one gets

$$\begin{cases} \mathbb{A} &= \mathbb{K}^{-1}(\mathbb{H} + \mathbf{1}_n \boldsymbol{\lambda}^\top) \\ \boldsymbol{\lambda}^\top &= \frac{1}{\delta}(\mathbf{1}_q^\top - \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbb{H}) \end{cases} \quad (7)$$

where $\mathbb{K} := \mathbb{E}[\mathbf{Y}^\top \mathbf{W} \mathbf{Y}]$, $\mathbf{h}(x^*) := \mathbb{E}[\mathbf{Y}^\top \mathbf{W} \mathbf{Y}(x^*)]$, and with scalar $\delta := \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbf{1}_n$.

For matrix expressions, $\boldsymbol{\lambda} := (\lambda(x_1^*), \dots, \lambda(x_q^*))^\top$, and $\mathbb{H} := [\mathbf{h}(x_1^*), \dots, \mathbf{h}(x_q^*)]$. 

If furthermore, for all $i = 1, \dots, p$, $x \in \mathcal{X}$, $\mathbb{E}[Y_i(x)] = \mu_i$, then $\mathbf{M}(x^)$ is unbiased.*

Optimal weights, constraints 1+2

Considered constraints

- **Constraint 1:** Weights sum to one, $\alpha^\top(x^*)\mathbf{1}_n = 1$, $x^* \in \mathcal{X}$.
- **Constraint 2:** Prescribed average \mathbf{m} of predicted values:

$$\mathbb{E}[\mathbf{M}(X^*)|\mathbb{Y}] = \mathbf{m}, \text{ with } X^* \text{ r.v. on } \{x_1^*, \dots, x_q^*\}, \text{ distribution } \pi.$$

Note: unlike usual kriging methods, weights **must** be calculated simultaneously.

Proposition (Joint Kriging weights under predicted values constraint)

The Joint Kriging weights minimizing the loss $\Delta(x^*)$ under the constraints 1+2 are:

$$\mathbb{A} = \mathbb{K}^{-1} \left(\mathbb{H} + \mathbf{1}_n \lambda^\top + \mathbb{Y}^\top \lambda' \pi^\top \right) \quad \text{NEW} \quad [\text{weights must be solved all at once}] \quad (8)$$

with Lagrange multipliers, provided that $\left(\frac{1}{\delta} \mathbf{u} \mathbf{u}^\top - \mathbb{Y} \mathbb{K}^{-1} \mathbb{Y}^\top \right)$ is invertible,

$$\lambda' = \gamma^{-1} \left(\frac{1}{\delta} \mathbf{u} \mathbf{u}^\top - \mathbb{Y} \mathbb{K}^{-1} \mathbb{Y}^\top \right)^{-1} \left(\mathbb{Y} \mathbb{K}^{-1} \mathbb{H} \pi + \frac{1}{\delta} \mathbf{u} \left(1 - \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbb{H} \pi \right) - \mathbf{m} \right)$$

$$\lambda = \delta^{-1} \left(\mathbf{1}_q - \mathbb{H}^\top \mathbb{K}^{-1} \mathbf{1}_n - \pi \lambda'^\top \mathbf{u} \right)$$

where $\mathbf{u} := \mathbb{Y} \mathbb{K}^{-1} \mathbf{1}_n$, $\gamma := \pi^\top \pi \in \mathbb{R}$ and $\delta := \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbf{1}_n \in \mathbb{R}$. $\pi = (P[X^* = x_i^*])_i$.

Affine extension (1/2)

External information source

Idea for prescribed \mathbf{m} : hidden external information (expert, other stat, etc.).
Let \mathbf{Z} be the random vector containing this source of information.

Affine predictor




The affine predictor is:

$$\mathbf{M}^+(x^*) := \alpha_0(x^*)\mathbf{Z} + \sum_{i=1}^n \alpha_i(x^*)\mathbf{Y}(x_i), \quad (9)$$

Given $\mathbf{Z} = \mathbf{m}$, a constant term is included in the sum, hence the name *affine prediction*.

Updated constraints

- **Constraint 1:** Weights sum to one, $\mathbf{1}_{n+1}^\top \boldsymbol{\alpha}^+(x^*) = 1$, $x^* \in \mathcal{X}$.
with $\boldsymbol{\alpha}^+ = (\alpha_0(x^*), \dots, \alpha_n(x^*))^\top$.
- **Constraint 2:** Prescribed average predicted values:

$$\mathbb{E} [\mathbf{M}^+(X^*) | \mathbf{Z} = \mathbf{m}, \mathbb{Y}] = \mathbf{m}, \text{ with } X^* \text{ r.v. on } \{x_1^*, \dots, x_q^*\} \text{ $$

Affine extension (2/2)

Updated optimal weights of the affine predictor can be derived easily:

Proposition (Affine version of predictors)

Assume that the following covariance vectors are given

$$\begin{cases} \mathbf{P}^\top & := & \mathbb{E} [\mathbf{Z}^\top \mathbb{W} \mathbf{Y}] - \mathbb{E} [\mathbf{Z}^\top] \mathbb{W} \mathbb{E} [\mathbf{Y}] \\ \mathbf{Q}^\top & := & \mathbb{E} [\mathbf{Z}^\top \mathbb{W} \mathbf{Y}^*] - \mathbb{E} [\mathbf{Z}^\top] \mathbb{W} \mathbb{E} [\mathbf{Y}^*] \\ \sigma_Z^2 & := & \mathbb{E} [\mathbf{Z}^\top \mathbb{W} \mathbf{Z}] - \mathbb{E} [\mathbf{Z}^\top] \mathbb{W} \mathbb{E} [\mathbf{Z}] \end{cases} \quad (10)$$

Then optimal weights of previous cases can be updated by replacing \mathbb{Y} , \mathbb{K} , \mathbb{H} by

$$\mathbb{Y}^+ = [\mathbf{m} \quad \mathbb{Y}], \quad \mathbb{K}^+ = \begin{bmatrix} \sigma_Z^2 & \mathbf{P}^\top \\ \mathbf{P} & \mathbb{K} \end{bmatrix}, \quad \mathbb{H}^+ = \begin{bmatrix} \mathbf{Q}^\top \\ \mathbb{H} \end{bmatrix}, \quad \text{NEW} \quad (11)$$

In practice one can set, e.g., $\sigma_Z^2 \ll \sigma^2$, and \mathbf{P} , \mathbf{Q} filled with zeros.

Kriging mean and variance

Recall the predictor shape, Kriging Mean, with optimal weights $\alpha(x^*)$

$$\mathbf{M}(x^*) := \mathbf{Y}\alpha(x^*) \quad (12)$$


Proposition (Joint Kriging variance with arbitrary weights)

Let $\alpha(x^*)$ be any vector of weights, possibly satisfying supplementary constraints. The associated Joint Kriging variance writes:

$$\Delta(x^*) = \alpha(x^*)^\top \mathbb{K}\alpha(x^*) - 2\alpha(x^*)^\top \mathbf{h}(x^*) + v(x^*), \quad (13)$$

or using a matrix expression, denoting $\mathbf{\Delta} := (\Delta(x_1^*), \dots, \Delta(x_q^*))^\top$, we get

$$\mathbf{\Delta} = \text{diag}[\mathbf{A}^\top \mathbb{K} \mathbf{A}] - 2 \text{diag}[\mathbf{A}^\top \mathbf{H}] + \text{diag}[\mathbb{K}^*],$$

where $\mathbb{K} := \mathbb{E}[\mathbf{Y}^\top \mathbb{W} \mathbf{Y}]$, $\mathbb{K}^* := \mathbb{E}[\mathbf{Y}^{*\top} \mathbb{W} \mathbf{Y}^*]$, $\mathbf{h}(x^*) := \mathbb{E}[\mathbf{Y}^\top \mathbb{W} \mathbf{Y}(x^*)]$ and $v(x^*) := \mathbb{E}[\mathbf{Y}(x^*)^\top \mathbb{W} \mathbf{Y}(x^*)]$ are given. Here $\mathbf{Y}^* := [\mathbf{Y}(x_1^*), \dots, \mathbf{Y}(x_q^*)]$. 

Here aggregated error, also variance sharing results for the error of each component.

Covariances (1/2)

Recall that all optimal weights rely on cross-moments matrices like $\mathbb{K} := \mathbb{E} [\mathbf{Y}^\top \mathbf{W} \mathbf{Y}]$.
One can easily replace these objects by “centered” ones:

Remark (Covariance matrices)

Assume that

(i) $\mathbb{E} [\mathbf{Y}(x)] = \boldsymbol{\mu}$ for all $x \in \mathcal{X}$,


(ii) either $\boldsymbol{\mu} = \mathbf{0}_p$ or weights sum to one $\boldsymbol{\alpha}(x^*)^\top \mathbf{1}_n = 1$.

Then the objects \mathbb{K} , \mathbb{H} , $\mathbf{h}(x^*)$, \mathbb{K}^* , $\mathbf{v}(x^*)$ can be replaced by the “centered” ones:

$$\left\{ \begin{array}{l} \tilde{\mathbb{K}} = \mathbb{E} [\mathbf{Y}^\top \mathbf{W} \mathbf{Y}] - \mathbb{E} [\mathbf{Y}^\top] \mathbf{W} \mathbb{E} [\mathbf{Y}] \\ \tilde{\mathbb{H}} = \mathbb{E} [\mathbf{Y}^\top \mathbf{W} \mathbf{Y}^*] - \mathbb{E} [\mathbf{Y}^\top] \mathbf{W} \mathbb{E} [\mathbf{Y}^*] \\ \tilde{\mathbf{h}}(x^*) = \mathbb{E} [\mathbf{Y}^\top \mathbf{W} \mathbf{Y}(x^*)] - \mathbb{E} [\mathbf{Y}^\top] \mathbf{W} \mathbb{E} [\mathbf{Y}(x^*)] \\ \tilde{\mathbb{K}}^* = \mathbb{E} [\mathbf{Y}^{*\top} \mathbf{W} \mathbf{Y}^*] - \mathbb{E} [\mathbf{Y}^{*\top}] \mathbf{W} \mathbb{E} [\mathbf{Y}^*] \\ \tilde{\mathbf{v}}(x^*) = \mathbb{E} [\mathbf{Y}(x^*)^\top \mathbf{W} \mathbf{Y}(x^*)] - \mathbb{E} [\mathbf{Y}(x^*)^\top] \mathbf{W} \mathbb{E} [\mathbf{Y}(x^*)] \end{array} \right. \quad (14)$$

everywhere in previous Propositions, without changing the optimal weights $\boldsymbol{\alpha}(x^*)$.

Covariances (2/2)

As $\Delta(x^*) \in \mathbb{R}$, the covariances rely on an implicit sum of components of $\mathbf{Y}(\cdot)$  :

$$k(x, x') := \mathbb{E} \left[\mathbf{Y}(x)^\top \mathbb{W} \mathbf{Y}(x') \right] - \mathbb{E} \left[\mathbf{Y}(x)^\top \right] \mathbb{W} \mathbb{E} \left[\mathbf{Y}(x') \right]$$

Remark (Using correlation functions)

Assume that there exists a positive definite matrix \mathbb{W} , such that the covariances depend only on some distance between x and x' :

$$k(x, x') = \sigma^2 r \left(\|x - x'\|_\theta \right) , \quad (15)$$

where $r(\cdot)$ is a unit correlation function and $\|x - x'\|_\theta^2 = \sum_{i=1}^d \left(\frac{x_i - x'_i}{\theta_i} \right)^2$ corresponds to a rescaled Euclidean norm. Then all components of the covariances matrices $\tilde{\mathbb{K}}$ and $\tilde{\mathbb{H}}$, etc. can be derived from the covariance function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$:

$$\begin{cases} \tilde{\mathbb{K}}_{ij} &= k(x_i, x_j) \\ \tilde{\mathbb{H}}_{ik} &= k(x_i, x_k^*) \\ \text{etc.} & \end{cases} \quad (16)$$

- Does not depend on \mathbb{W} any more: with this assumption, **no need to estimate \mathbb{W}** .
- This simplifies **a lot** the hyperparameters estimation.

Constrained Classification

Application to constraint classification (1/3)

Assumption on observations





- **Label binarization.** Each class label $\ell \in \{1, \dots, p\}$ can be converted into a vector of indicator functions (even for non ordinal classes):

$$\mathbf{Y} := \left(\mathbb{1}_{\{j=\ell\}} \right)_{j=1, \dots, p}.$$

- **Observation of membership degrees.** More generally, each observation $\mathbf{Y}(x_i)$ consists in a distribution of possible labels, where degrees sum to one

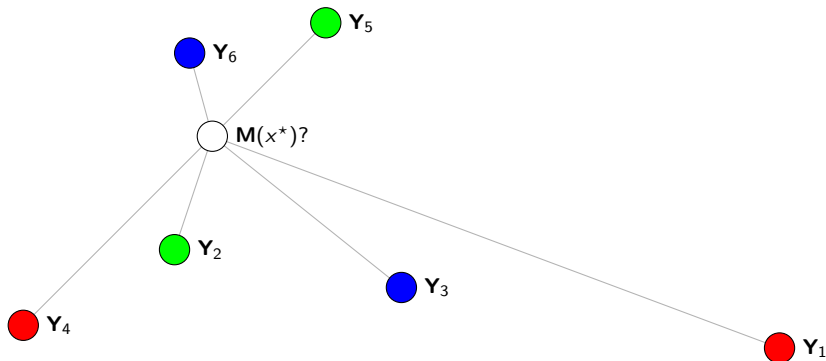
“Constraint 3”: $\mathbf{1}_p^\top \mathbf{Y}(x_i) = 1, \quad i = 1, \dots, n.$

Non ordinal example of observations, $p = 3$ classes: {red, green, blue}

 $\mathbf{Y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix},$
 $\mathbf{Y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix},$
 $\mathbf{Y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$
 and even  $\mathbf{Y} = \begin{bmatrix} 1/2 \\ 0 \\ 1/2 \end{bmatrix}$

only requirement: components sum to one.

Application to constraint classification (2/3)

Illustration of the fuzzy classification

Predicted membership degrees at point x^* :

$$M(x^*) = \underbrace{\alpha_1(x^*)}_{\in \mathbb{R}} Y_1 + \dots + \underbrace{\alpha_6(x^*)}_{\in \mathbb{R}} Y_6 = 2\% \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \dots + 40\% \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 10\% \\ 35\% \\ 55\% \end{bmatrix}$$

... weights $\alpha_i(x^*)$ are obtained by Joint Kriging formulas, under chosen constraints.

Application to constraint classification (3/3)

Constrained Classification: apply Joint Kriging model to membership degrees  :

Remark (Constraints impact)

Recall all constraints

- *constraint 1: weights sum to one,*
- *constraint 2: prescribed average of predictions,*
- *constraint 3: observations are membership degrees, $\mathbf{1}_p^\top \mathbb{Y} = \mathbf{1}_q^\top$.*

Then with constrained Joint Kriging model:

- *Predicted membership degrees are summing to one:*

$$\text{Constraints 1+3} \implies \mathbf{1}_p^\top \mathbf{M}(x^*) = 1, \quad x^* \in \mathcal{X}$$

- *Average class percentages over prediction points can be chosen:*

$$\text{Constraints 2+3} \implies \mathbb{E}[\mathbf{M}(X^*)|\mathbb{Y}] = \mathbf{m}, \quad \text{with } \mathbf{1}_p^\top \mathbf{m} = 1$$

\mathbf{m} is the prescribed average of each class, and X^ a rv over all prediction points.*

Numerical Illustrations

available notebooks



Available notebooks

All illustrations are generated with notebooks that are available at <https://gitlab.emse.fr/marc.grossouvre/jointkrigingsupplementary/>

In **modifiable executable** format `.Rmd` and in **executed directly readable** `.html` format.

We did our best to make results **fully reproducible**, and figures settings easy to retrieve.

A simple sinus function, 1D input, 1D output

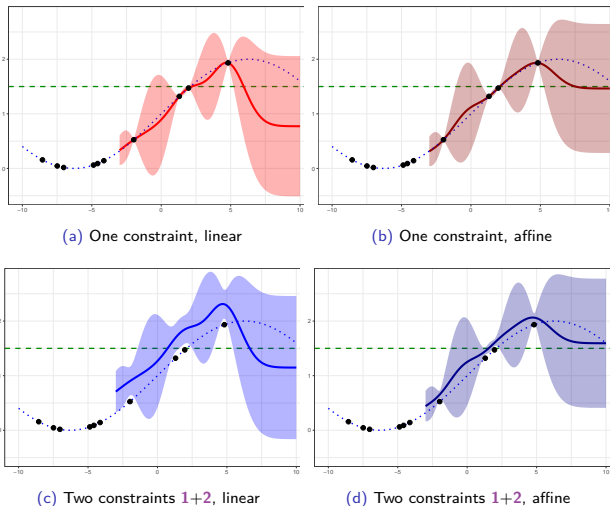


Figure: Joint Kriging Prediction. prescribed value $m = 1.5$ (horizontal dashed line). Observations are black dots, the thin dotted blue line is the underlying function.

Time-series Air Quality data, 1D input, 4D output, Constraint 1

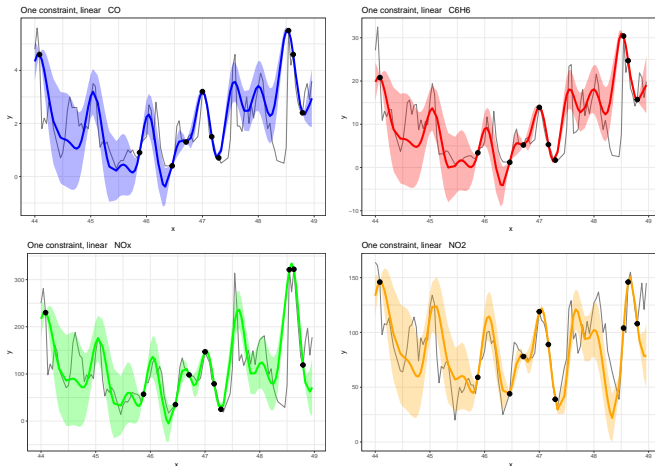


Figure: Joint Kriging interpolation with Constraint 1. data points (black dots). Top: CO, C6H6, bottom: NOx, NO2. Predictions thick solid lines, true values are in thin black solid lines.

Time-series Air Quality data, 1D input, 4D output, chosen covariance

Despite $p = 4$ outputs, one needs a **single** covariance structure (which involves linear combinations of components).

Multiply a **periodic kernel with period of one day**, and Matérn 3/2 kernel, parameter θ :

$$k(x, x') = \sigma^2 \exp(-\sin^2(\pi|x - x'|)) \left(1 + \frac{|x - x'|}{\theta}\right) \exp\left(-\frac{|x - x'|}{\theta}\right). \quad (17)$$

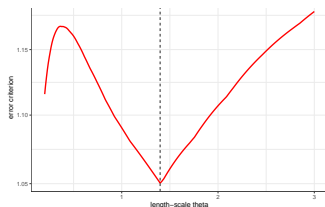


Figure: Optimization of the *single* correlation hyperparameter θ for the four selected pollutants, data extracted from Air quality data set.

Caution: depends quite heavily on the chosen observation locations, sometimes the error function is monotonic! easier to control with very few hyperparameters!

Time-serie Air Quality data, 1D input, 4D output, Constraint 1+2 adverse

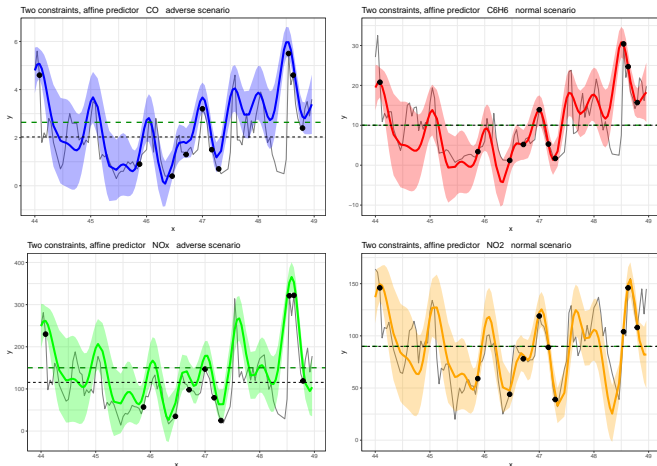


Figure: Adverse scenarios: constraints 1+2 and affine predictor. Left panels: adverse scenarios, average 130% of the true average, right panels: regular scenarios.

Time-series Air Quality data, 1D input, 4D output

Difficult problem

predicting four quite erratic time series from 10 observations

Small experiment conclusion

- can handle complex multi-valued data
- satisfies all chosen constraints
- very simple, closed-form formulas
- has a limited number of hyperparameters
- allows adverse scenarios
- performs reasonably well

but no benchmark yet... done in the next experiment.

Quake classification problem, input 3D, output "2D"

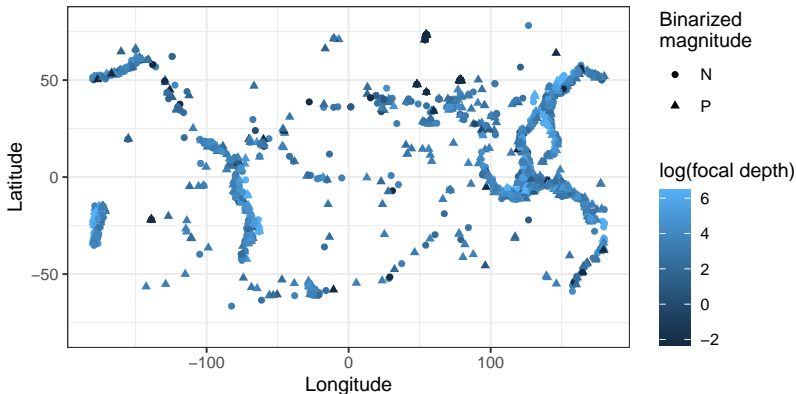


Figure: Earthquakes observations. An earthquake is a point with coordinates latitude, longitude and focal depth (given by the color). Triangles: earthquakes which magnitude is above average. Circles: below average.

data available at www.openml.org/search?type=data&id=772.

all notebooks: <https://gitlab.emse.fr/marc.grossouvre/jointkrigingsupplementary/>

Quake classification problem - covariance

Chosen covariance function

A **single** covariance function (only one tested!):

$$k(x, x') = \sigma^2 \exp \left(-2 \frac{\sin^2((x_1 - x'_1)/2)}{\theta_1^2} - 2 \frac{\sin^2((x_2 - x'_2)/2)}{\theta_2^2} \right) \exp \left(-2 \frac{(x_3 - x'_3)^2}{\theta_3^2} \right)$$

periodicity of longitude, latitude, not focal depth. Small nugget (rounded magnitudes).

Parameter estimation

The hyperparameters estimation has been treated separately on other train/test splits to avoid overfitting the data (using 10 fold cross-validation).

Resulting values for θ are 2.3 for latitude, 0.9 for longitude and 196.8 for focal depth.

Quake classification problem - prediction constraints 1+2

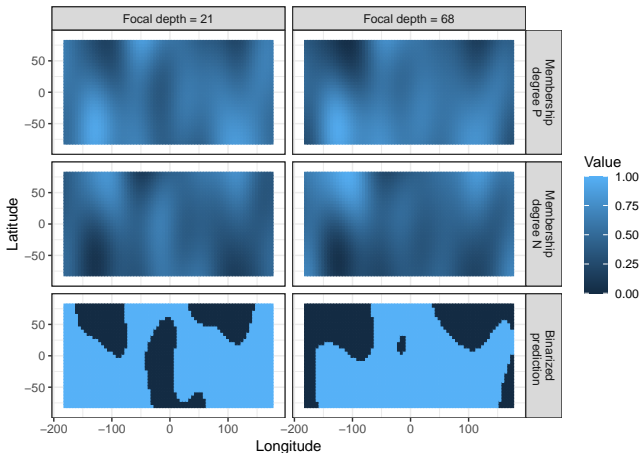


Figure: Joint Kriging with two constraints 1+2. Top: membership degree of “P: magnitude is above average”, bottom: membership degree of “N: magnitude is below average”, binarized prediction (1 if membership degree of P is greater than 0.5). Left: 21km focal depth (=Q1). Right: 68km (=Q3).

Quake classification problem - benchmark

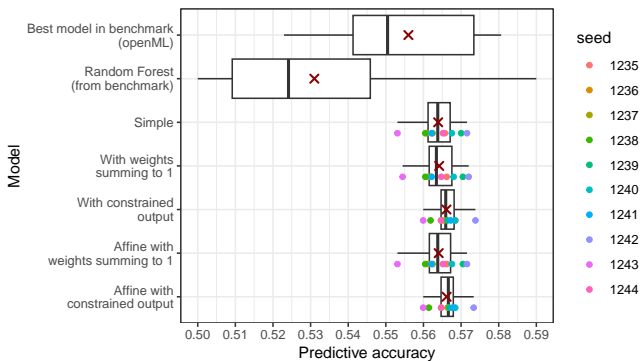


Figure: Performances for 10 runs. Top 2: best OpenML model (kernel logistic regression) and OpenML Random Forest. Bottom 5: Joint Kriging models. Dark red cross: average predictive accuracy, **the higher the better**.

Predictive accuracy for Joint Kriging: 0.5661 ± 0.0038 . For best OpenML: 0.556 ± 0.018 .
 OpenML benchmark: 69 models, www.openml.org/search?type=task&id=4516.

Quake classification problem - adverse scenario

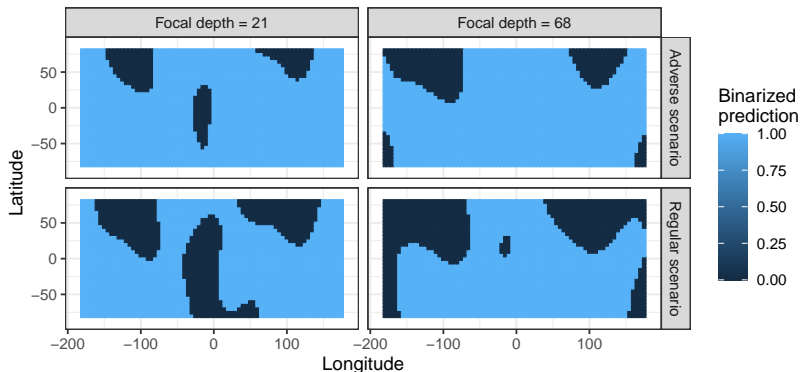


Figure: Adverse scenario, constraints 1+2. Top: adverse scenario, first class output average constrained to be 65%. Bottom: regular scenario, output average constrained to 55.5%. Left: 21km focal depth. Right: 68km focal depth.

Quake classification problem - four classes

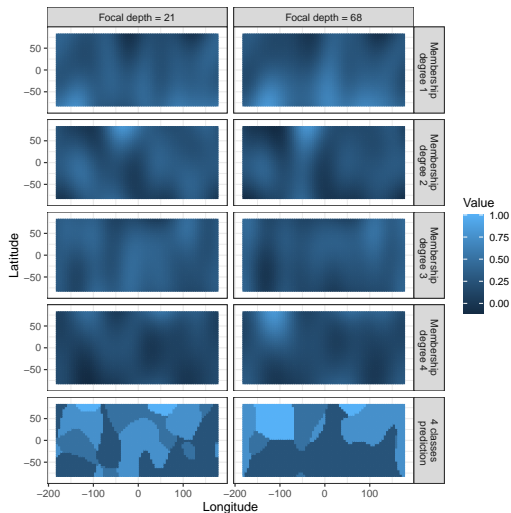


Figure: Affine Joint Kriging with constraints 1+2, magnitude thresholds 5.85, 5.95, 6.15. bottom: class of greatest membership degree. Remark longitude circular coherence.

Conclusion

what is done

- **Multi-output** Kriging model, not necessarily Gaussian.
- **Specific simplification**: weights apply jointly to all outputs.
- **Specific constraints**, especially on predicted values.
- **Specific affine model**.

Pros

- **Simple**: computable in closed-form, drastically reduces hyperparameters number.
- **Useful**: interpretable, can interpolate data, uncertainty measurements, specific covariances (e.g. periodicity). Constraints allows for external information, expert judgments, adverse modelling, or homogenization needs such as fairness constraints. For fuzzy classification, prescribed class percentages.
- **Performant**: competes with state-of-the-art algorithms on an open benchmark.

Cons and perspectives

- **Simplified**: possible limitations for different regularities of outputs. Introduce more complex covariances, model with higher complexity...
- **Needs hyperparameters optimization**: specific estimation procedures...
- **Non-convex**: possible membership degrees outside $[0, 1]$. Convex constraints...
- **Broken continuous interpolation property** with Constraint 2. Modify predictor...

Thank you for your attention !

- Questions?
- Details and proofs in the preprint <https://hal.science/hal-04208454>.
- Do not hesitate to send comments or references!

References I

- [1] Agarwal, G., Sun, Y., and Wang, H. J. (2021). Copula-based multiple indicator kriging for non-gaussian random fields. *Spatial Statistics*, 44:100524.
- [2] Benatti, K. A., Pedroso, L. G., and Ribeiro, A. A. (2022). Theoretical analysis of classic and capacity constrained fuzzy clustering. *Information Sciences*, 616:127–140.
- [3] Bradley, P. S., Bennett, K. P., and Demiriz, A. (2000). Constrained k-means clustering. *Microsoft Research, Redmond*, 20(0):0.
- [4] Chiang, J.-L., Liou, J.-J., Wei, C., and Cheng, K.-S. (2013). A feature-space indicator kriging approach for remote sensing image classification. *IEEE transactions on geoscience and remote sensing*, 52(7):4046–4055.
- [5] Dahl, A. and Bonilla, E. V. (2019). Grouped gaussian processes for solar power prediction. *Machine Learning*, 108(8-9):1287–1306.
- [6] Furrer, R. and Genton, M. G. (2011). Aggregation-cokriging for highly multivariate spatial data. *Biometrika*, 98(3):615–631.
- [7] Ganganath, N., Cheng, C.-T., and Tse, C. K. (2014). Data clustering with cluster size constraints using a modified k-means algorithm. In *2014 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery*, pages 158–161.

References II

- [8] Goovaerts, P. (1998). Ordinary cokriging revisited. *Mathematical Geology*, 30:21–42.
- [9] Goovaerts, P. (2009). Auto-ik: A 2d indicator kriging program for the automated non-parametric modeling of local uncertainty in earth sciences. *Computers & Geosciences*, 35(6):1255–1270.
- [10] Gordon, A. (1996). A survey of constrained classification. *Computational Statistics & Data Analysis*, 21(1):17–29.
- [11] Höppner, F. and Klawonn, F. (2008). Clustering with size constraints. In Jain, L. C., Sato-Ilic, M., Virvou, M., Tsihrintzis, G. A., Balas, V. E., and Abeynayake, C., editors, *Computational Intelligence Paradigms: Innovative Applications*, pages 167–180. Springer Berlin Heidelberg, Berlin, Heidelberg.
- [12] Journel, A. G. (1983). Nonparametric estimation of spatial distributions. *Journal of the International Association for Mathematical Geology*, 15:445–468.
- [13] Meer, F. V. D. (1996). Classification of remotely-sensed imagery using an indicator kriging approach: application to the problem of calcite-dolomite mineral mapping. *International Journal of Remote Sensing*, 17(6):1233–1249.
- [14] Panos, A., Dellaportas, P., and Titsias, M. K. (2021). Large scale multi-label learning using gaussian processes. *Machine Learning*, 110:965–987.
- [15] Rasmussen, C. E., Williams, C. K., et al. (2006). *Gaussian processes for machine learning*, volume 1. Springer.

References III

- [16] Ver Hoef, J. M. and Cressie, N. (1993). Multivariable spatial prediction. *Mathematical Geology*, 25:219–240.
- [17] Williams, C. K. and Barber, D. (1998). Bayesian classification with gaussian processes. *IEEE Transactions on pattern analysis and machine intelligence*, 20(12):1342–1351.
- [18] Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. (2019). Fairness constraints: A flexible approach for fair classification. *The Journal of Machine Learning Research*, 20(1):2737–2778.

Notations I

locations

\mathcal{X} set of locations (inputs/design points).

n, q number of observed locations, of prediction locations.

x any location. x_1, \dots, x_n are all observed locations.

x^* any prediction location. x_1^*, \dots, x_q^* are all prediction locations.

X^* a random variable over prediction locations.

$\pi = (\pi_{x_1^*}, \dots, \pi_{x_q^*})$ the $q \times 1$ distribution of X^* over prediction locations.

$\gamma = \pi^\top \pi$ an intermediate real value used in calculations.

targets

p number of targets (i.e. number of outputs).

$\mathbf{Y}(x)$ the $p \times 1$ vector of targets at location x .

$\mu = \mathbb{E}[\mathbf{Y}(x)]$ the $p \times 1$ mean of $\mathbf{Y}(x)$, when constant over x .

$\mathbb{Y} = [\mathbf{Y}(x_1), \dots, \mathbf{Y}(x_n)]$ all the $p \times n$ values of observed targets.

$\mathbb{Y}^* = [\mathbf{Y}(x_1^*), \dots, \mathbf{Y}(x_q^*)]$ all $p \times q$ unknown targets at prediction locations.

prediction

$\mathbf{M}(x^*)$ a $p \times 1$ predictor of $\mathbf{Y}(x)$

$\mathbb{M} = [\mathbf{M}(x_1^*), \dots, \mathbf{M}(x_q^*)]$ the $p \times q$ matrix of all predictions.

$\alpha(x^*)$ the $n \times 1$ linear weights for the prediction in x^* .

$\mathbb{A} = [\alpha(x_1^*), \dots, \alpha(x_q^*)]$ the $n \times q$ matrix of weights for all predictions.

\mathbf{m} a given constant $p \times 1$ vector of prescribed mean predicted values.

Notations II

$\Delta(x^*)$ loss to be minimized for finding $\mathbf{M}(x^*)$.

λ a $q \times 1$ vector of Lagrange multipliers (relative to sum of weights)

λ' a $p \times 1$ vector of Lagrange multipliers (relative to predicted values)

$\mathbf{u} = \mathbf{Y}\mathbf{K}^{-1}\mathbf{1}_n$ an intermediate $p \times 1$ vector in calculations.

\mathbf{Z} an additional $p \times 1$ factor for affine predictions.

covariances

$k(.,.)$ a covariance function.

\mathbb{W} a given symmetric positive definite matrix for computing norms.

$\mathbf{h}(x^*) = \mathbf{E} [\mathbf{Y}^\top \mathbb{W} \mathbf{Y}(x^*)]$ a $n \times 1$ covariance vector.

$\mathbb{H} = (\mathbf{h}(x_1^*), \dots, \mathbf{h}(x_q^*))$ a $n \times q$ covariance matrix.

$\mathbb{K} = \mathbf{E} [\mathbf{Y}^\top \mathbb{W} \mathbf{Y}]$ a $n \times n$ covariance matrix.

$\tilde{\mathbb{K}}, \tilde{\mathbf{h}}(x^*), \tilde{\mathbb{H}}$ other covariances using centred expressions.

$\delta = \mathbf{1}_n^\top \mathbb{K}^{-1} \mathbf{1}_n$ an intermediate real value in calculations.

\mathbf{P} additional $n \times 1$ covariance vector between \mathbf{Z} and $\mathbf{Y}(x_i)$

\mathbf{Q} additional $q \times 1$ covariance vector between \mathbf{Z} and $\mathbf{Y}(x_j^*)$

miscellaneous

\mathbf{v} a generic vector for defining norm or checking psd characteristic.

$\mathbf{1}_n, \mathbf{1}_p, \mathbf{1}_q$ a vector of ones of size n, p, q respectively.

$\mathbf{0}_n, \mathbf{0}_p, \mathbf{0}_q$ a vector of zeros of size n, p, q respectively.