

# Formation CIROQUO

---

Tanguy APPRIOU <sup>(1), (2)</sup>

(1) STELLANTIS

(2) École des Mines de Saint-Etienne, LIMOS

18 novembre 2024





### R Package “combinedkriging”

Available on Github : <https://github.com/TAppriou/combinedkriging>

Notebook also on Github :

Or with the link, then make a copy :

<https://colab.research.google.com/drive/1jMLOrt3PTepAgABOJGTz0jiu5-fueJXb?usp=sharing>

- Package to build a combination of Kriging models with fixed length-scales.

→ Use for high-dimensional Bayesian optimization.



→ We propose a model which is a combination of Kriging models with random length-scales

$$M_{tot}(\mathbf{x}) = \sum_{i=1}^p w_i(\mathbf{x}) M_i(\mathbf{x}),$$

with  $M_i(\mathbf{x}) = \mu_i + k_{\theta_i}(\mathbf{x}, \mathbf{X}) K_{\theta_i}^{-1}(\mathbf{Y} - \mu_i)$  Kriging model with fixed length-scale vector  $\theta_i$ .

### 1) Choice of the weights

## CHOICE OF THE WEIGHTS

- Weights based on the sub-models variance:  $\hat{s}_i^2(\mathbf{x}) = k_{\theta_i}(\mathbf{x}, \mathbf{x}) - k_{\theta_i}(\mathbf{x}, \mathbf{X})\mathbf{K}_{\theta_i}(\mathbf{X}, \mathbf{X})^{-1}k_{\theta_i}(\mathbf{X}, \mathbf{x})$ 
  - Product of Experts (PoE): correspond to the weights of the best linear combination for independent sub-models.

$$w_{PoE_i}(\mathbf{x}) = \frac{\hat{s}_i^{-2}(\mathbf{x})}{\sum_{j=1}^p \hat{s}_j^{-2}(\mathbf{x})}$$

→ Do not depend on the observations.

→ Gives more weight to models with higher length-scales.

- Generalized PoE (gPoE): corrects PoE by weighting the contribution of each sub-models.

$$w_{gPoE_i}(\mathbf{x}) = \frac{\beta_i \hat{s}_i^{-2}(\mathbf{x})}{\sum_{j=1}^p \beta_j \hat{s}_j^{-2}(\mathbf{x})}, \text{ the inner weights } \boldsymbol{\beta} \text{ can be obtained by LOOCV for example :}$$

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} e_{LOOCV} \left( \sum_{i=1}^p w_{gPoE_i}(\boldsymbol{\beta}) M_i \right), \quad \text{subject to } \sum_{i=1}^p w_{gPoE_i}(\boldsymbol{\beta}) = 1.$$

→ No analytical expression.

## CHOICE OF THE WEIGHTS

- Mixture of experts (MoE): weights based on the sub-models likelihoods.

$$w_{MoE_i} = \frac{\mathcal{L}(\boldsymbol{\theta}_i)}{\sum_{j=1}^p \mathcal{L}(\boldsymbol{\theta}_j)}, \quad \text{where } \mathcal{L}(\boldsymbol{\theta}) = -\frac{1}{2} \mathbf{Y}^T \mathbf{K}_{\boldsymbol{\theta}}^{-1} \mathbf{Y} - \frac{1}{2} \log |\mathbf{K}_{\boldsymbol{\theta}}| - \frac{n}{2} \log(2\pi)$$

→ The likelihood is not always a good measure of model accuracy when few observations are available.

→ The likelihood of sub-models often differ by several order of magnitudes, thus this method often select a few sub-models instead of doing a combination.

## CHOICE OF THE WEIGHTS

- Weights based on LOOCV: 
$$e_{LOOCV}(M_{tot}) = \frac{1}{n} \sum_{k=1}^n \left( \sum_{i=1}^p w_i M_{i-k}(x_k) - y(x_k) \right)^2 = \mathbf{w}^T \mathbf{C} \mathbf{w}.$$

→ The components of the matrix  $\mathbf{C}$  are :  $c_{ij} = \frac{1}{N} e_{CV_i}^T e_{CV_j}$ , with  $e_i^{(k)} = [K_{\theta_i}^{-1}(Y - \mu_i)]_k / [K_{\theta_i}^{-1}]_{k,k}$ ,  $k = 1, \dots, n$

- Normal LOOCV:

$$\mathbf{w}_{LOOCV} = \arg \min_{\mathbf{w}} \mathbf{w}^T \mathbf{C} \mathbf{w}, \quad \text{subject to } \mathbf{1}^T \mathbf{w} = 1 \quad \Rightarrow \mathbf{w}_{LOOCV} = \frac{\mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}$$

→ Weights can be negative or more than 1.

- Diagonal LOOCV: to enforce  $w_i \in [0,1]$ , we keep only the diagonal terms of  $\mathbf{C}$  :

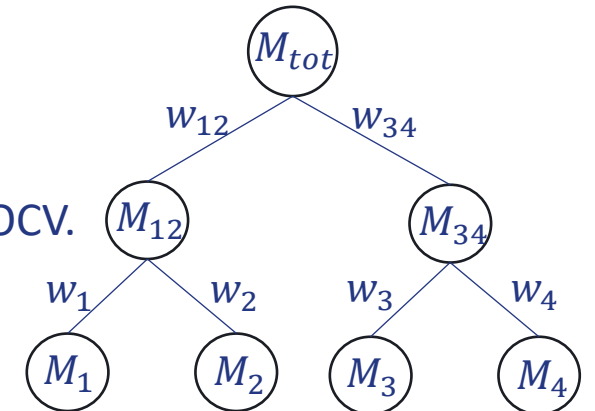
$$\mathbf{w}_{LOOCV_{diag}} = \frac{\mathbf{C}_{diag}^{-1} \mathbf{1}}{\mathbf{1}^T \mathbf{C}_{diag}^{-1} \mathbf{1}} \quad \Rightarrow w_{LOOCV_{diag_i}} = \frac{e_{LOOCV_i}^{-1}}{\sum_{l=1}^P e_{LOOCV_l}^{-1}}$$

→ Tends to give weights to all models (close to average of all models)

- Binary LOOCV: to enforce  $w_i \in [0,1]$ , we combine the sub-models two by two using LOOCV.

→ Number of sub-models limited to a power of 2.

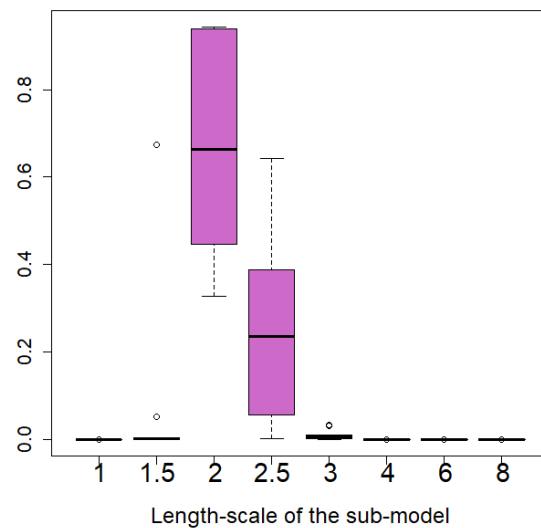
→ Can implement sparsity.



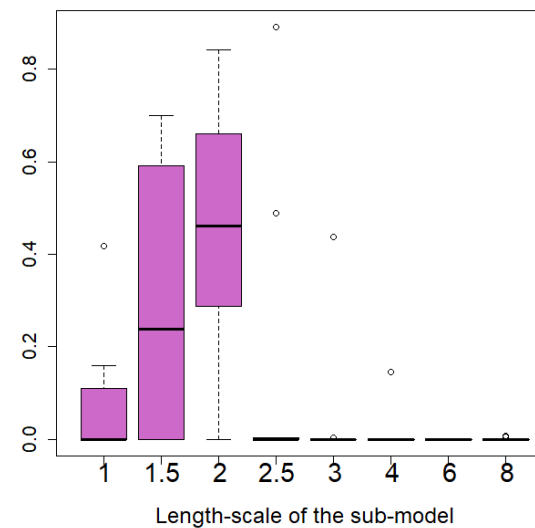
# CHOICE OF THE WEIGHTS



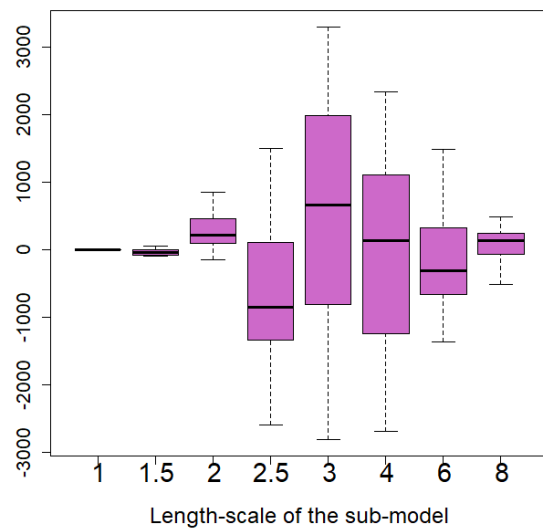
MoE weights for isotropic models



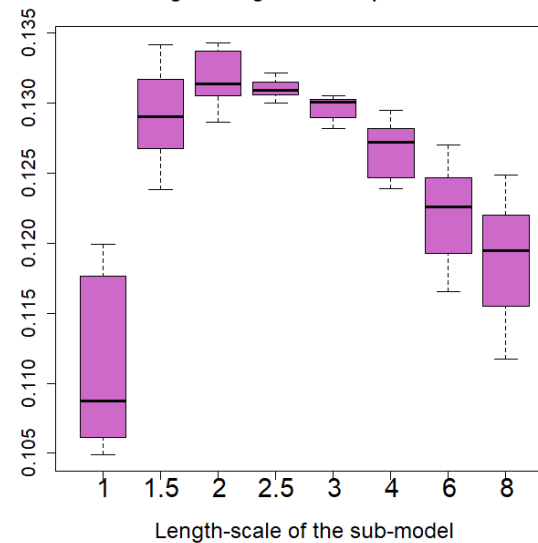
gPoE weights for isotropic models



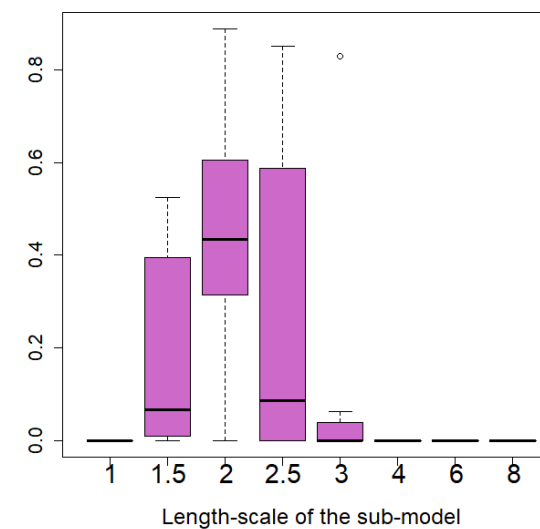
LOOCV weights for isotropic models



diagLOO weights for isotropic models



binLOO weights for isotropic models





→ We propose a model which is a combination of Kriging models with random length-scales

$$M_{tot}(\mathbf{x}) = \sum_{i=1}^p w_i(\mathbf{x}) M_i(\mathbf{x}),$$

with  $M_i(\mathbf{x}) = \mu_i + k_{\theta_i}(\mathbf{x}, \mathbf{X}) K_{\theta_i}^{-1}(\mathbf{Y} - \mu_i)$  Kriging model with fixed length-scale vector  $\theta_i$ .

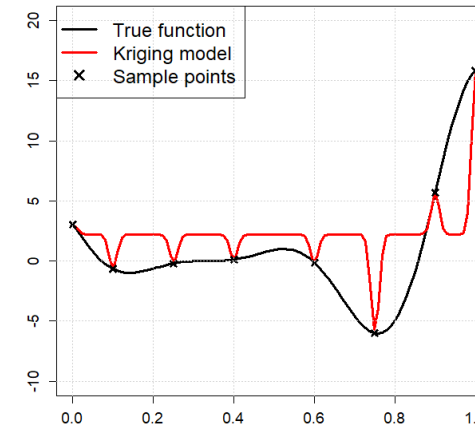
### 2) Choice of the sub-model length-scales



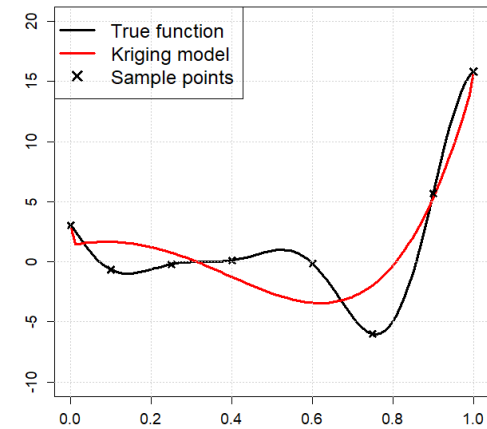
## SAMPLING THE LENGTH-SCALES

- We want to sample the length-scales in a **range of appropriate values** to avoid degenerate cases.
- For too small values:  $k_\theta(x_i, x_j) \rightarrow 0$  for all  $i \neq j$ , and  $\mathbf{K}_\theta \rightarrow \sigma^2 \mathbf{I}_n$ .
- For too large values:  $k_\theta(x_i, x_j) \rightarrow 1$ , and  $\mathbf{K}_\theta \rightarrow \mathbf{1}_{n \times n}$ .

Small  $\theta$



Large  $\theta$  (+nugget)



→ We need to define an interval  $[\theta_{min}^{(\ell)}, \theta_{max}^{(\ell)}]$  for sampling the length-scales.



- Assume the random vector of design points is  $\mathbf{X} = (X^{(1)}, \dots, X^{(d)})$  with i.i.d components.

We denote  $\sigma^2 = \text{Var}(X^{(i)})$ , and  $\kappa = E \left[ \left( \frac{X^{(i)} - \mu}{\sigma} \right)^4 \right]$ .

We can obtain the distribution of the distance between two random points  $\mathbf{X}$  and  $\mathbf{X}'$  :

$$D^2 = \sum_{i=1}^d (X^{(i)} - X'^{(i)})^2 \sim \mathcal{N}(2d\sigma^2, 2d\sigma^4(\kappa + 1))$$

- Typical distances in the DoE can be obtained by taking the root of a 95% Gaussian confidence interval :

$$[r_{min}, r_{max}] = \left[ \sigma \sqrt{2d - 1,96\sqrt{2(\kappa + 1)d}}, \sigma \sqrt{2d + 1,96\sqrt{2(\kappa + 1)d}} \right].$$

- If the standard deviation differs on every dimension  $\ell$  :

$$[r_{min}^{(\ell)}, r_{max}^{(\ell)}] = \sigma^{(\ell)} [r_{min}, r_{max}] = \left[ \sigma^{(\ell)} \sqrt{2d - 1,96\sqrt{2(\kappa + 1)d}}, \sigma^{(\ell)} \sqrt{2d + 1,96\sqrt{2(\kappa + 1)d}} \right]$$

- For a fixed distance  $r$ , the impact of small variations of the length-scale on the correlation can be measure by the index :

$$I\left(\frac{r}{\theta}\right) = \left| \frac{\frac{\partial}{\partial \theta} k\left(\frac{r}{\theta}\right)}{\max_{\theta} \frac{\partial}{\partial \theta} k\left(\frac{r}{\theta}\right)} \right|$$

- For a given distance, a length-scale is considered influential if the index is more than a threshold for the typical distances :

$$\Theta_{adm}(r) = \left\{ \theta : I^{(\ell)}\left(\frac{r}{\theta}\right) \geq \delta \right\}.$$

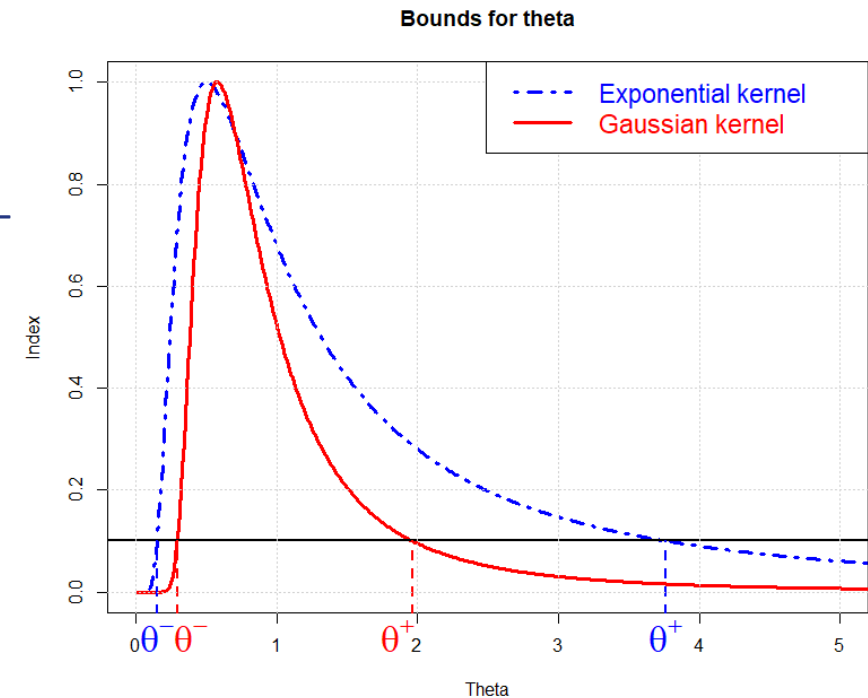
- Note that multiplying by a scale factor  $\alpha$  changes the set of admissible length-scales by the same factor :

$$\Theta_{adm}(\alpha r) = \alpha \Theta_{adm}(r).$$

→ We only have to solve for  $r = 1$  in  $\theta$  :

$$I^{(\ell)}\left(\frac{1}{\theta}\right) = \delta.$$

- We denote  $\theta^-(k)$  and  $\theta^+(k)$  the smallest and largest roots of  $I^{(\ell)}\left(\frac{1}{\theta}\right) = \delta$ .



## SAMPLING THE LENGTH-SCALES

- Putting both factors together :

For typical values of the inter-point distance  $r \in [r_{min}^{(\ell)}, r_{max}^{(\ell)}]$ , the length-scales bounds are chosen as :

$$\theta_{min}^{(\ell)} = \inf_{r \in [r_{min}^{(\ell)}, r_{max}^{(\ell)}]} \Theta_{adm}(r) = r_{min}^{(\ell)} \theta^-(k), \quad \theta_{max}^{(\ell)} = \sup_{r \in [r_{min}^{(\ell)}, r_{max}^{(\ell)}]} \Theta_{adm}(r) = r_{max}^{(\ell)} \theta^+(k).$$

Finally, we get :

$$\theta_{min}^{(\ell)} = \underbrace{\sigma^{(\ell)} r_{min}}_{\text{Influence of the design}} \underbrace{\theta^-(k)}_{\text{Influence of covariance family}} \quad \text{and} \quad \theta_{max}^{(\ell)} = \underbrace{\sigma^{(\ell)} r_{max}}_{\text{Influence of the design}} \underbrace{\theta^+(k)}_{\text{Influence of covariance family}}.$$

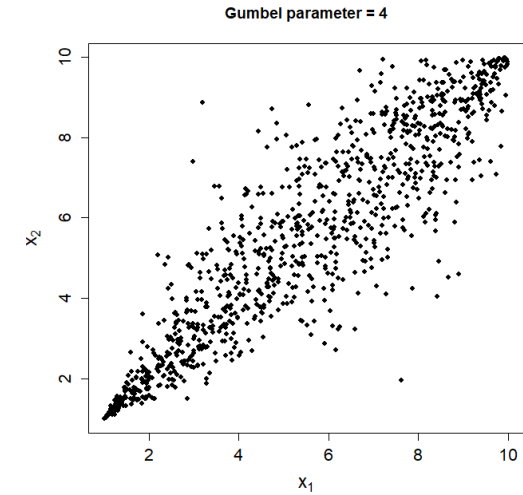
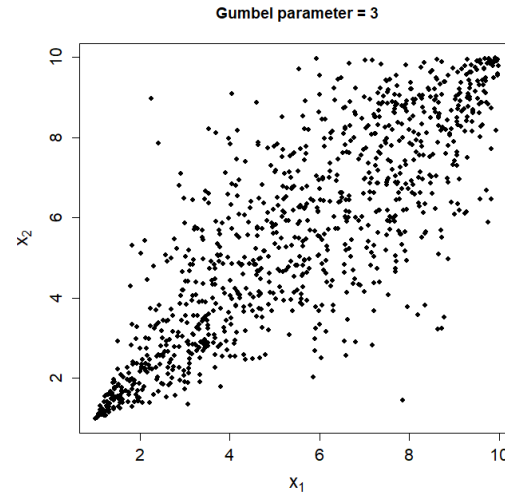
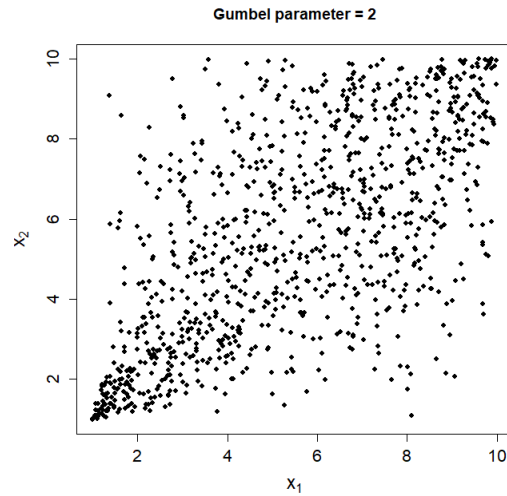
d	Kernel $k$	Design influence			Kernel influence		Resulting bounds	
		$\sigma^{(\ell)}$	$r_{min}$	$r_{max}$	$\theta^-(k)$	$\theta^+(k)$	$\theta_{min}^{(\ell)}$	$\theta_{max}^{(\ell)}$
10	Exponential	$\frac{1}{\sqrt{12}}$	2.31	5.89	0.15	3.76	0.10	6.39
	Matérn 3/2				0.21	2.74	0.14	4.66
	Matérn 5/2				0.23	2.44	0.15	4.15
	Gaussian				0.29	1.96	0.19	3.33
50	Exponential	$\frac{1}{\sqrt{12}}$	8.20	11.5	0.15	3.76	0.36	12.5
	Matérn 3/2				0.21	2.74	0.50	9.10
	Matérn 5/2				0.23	2.44	0.54	8.10
	Gaussian				0.29	1.96	0.69	6.51
$d \rightarrow \infty$	Exponential	$\frac{1}{\sqrt{12}}$	$\sqrt{2d}$	$\sqrt{2d}$	0.15	3.76	$0.061\sqrt{d}$	$1.54\sqrt{d}$
	Matérn 3/2				0.21	2.74	$0.086\sqrt{d}$	$1.12\sqrt{d}$
	Matérn 5/2				0.23	2.44	$0.094\sqrt{d}$	$1.00\sqrt{d}$
	Gaussian				0.29	1.96	$0.12\sqrt{d}$	$0.80\sqrt{d}$

Example of values for a uniform design plan ( $\kappa = 9/5$ ), a standard deviation  $\sigma^{(\ell)} = 1/\sqrt{12}$ , and a threshold  $\delta = 1/10$ .

## SAMPLING THE LENGTH-SCALES

- Uniform sampling:  $\theta^{(\ell)} \sim \mathcal{U}[\theta_{\min}^{(\ell)}, \theta_{\max}^{(\ell)}], \ell = 1, \dots, d.$

- Sampling with copulas:



- Sampling based on the entropy of the correlations

- Recall that  $D^2$  the random square distance between two independent points  $\mathbf{X}$  and  $\mathbf{X}'$  of the design is :

$$D^2 = \sum_{k=1}^d (X_k - X'_k)^2 \sim \mathcal{N} \left( 2d\sigma_X^2, 2d\sigma_X^4(\kappa_X + 1) \right).$$

- For a Gaussian correlation (for other kernels, we can use numerical approximations):

$$R_\theta = e^{-\frac{1D^2}{2\theta^2}} \sim \log \mathcal{N} \left( \frac{-\sigma_X^2}{\theta^2} d, \frac{\sigma_X^4}{2\theta^4} (\kappa_X + 1)d \right).$$

- We can finally obtain the entropy of the correlation:

$$H(R_\theta) = \mathbb{E}(-\log f_{R_\theta}(R_\theta)) = -\frac{\sigma_X^2}{\theta^2} d + \frac{1}{2} \ln \left( \frac{\sigma_X^4}{2\theta^4} d(\kappa_X + 1)2\pi \right) + \frac{1}{2}.$$

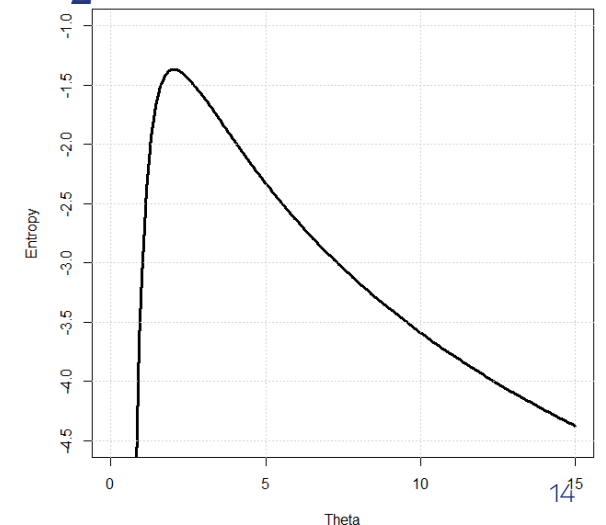
- When sampling the length-scales, we want to favor  $\theta$  corresponding **to high entropy values, which result in a high variability in the correlation.**

→ We will sample the length-scales using a positive transformation of the entropy:

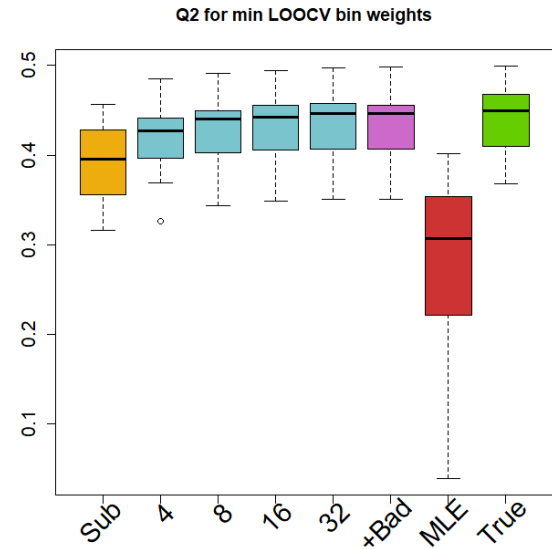
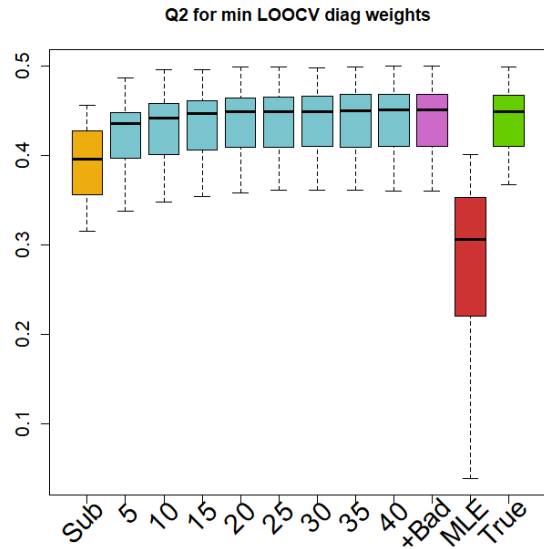
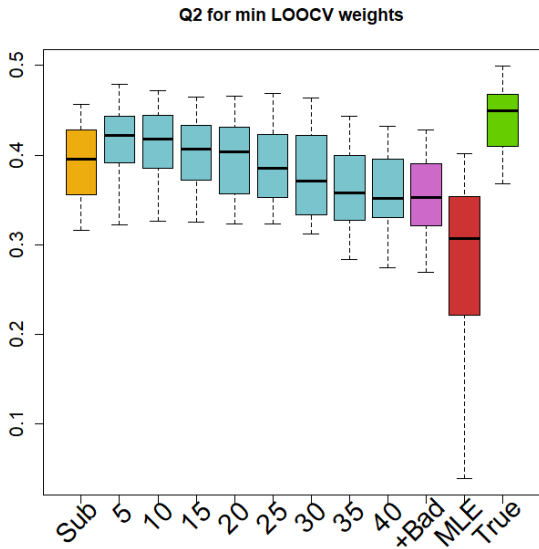
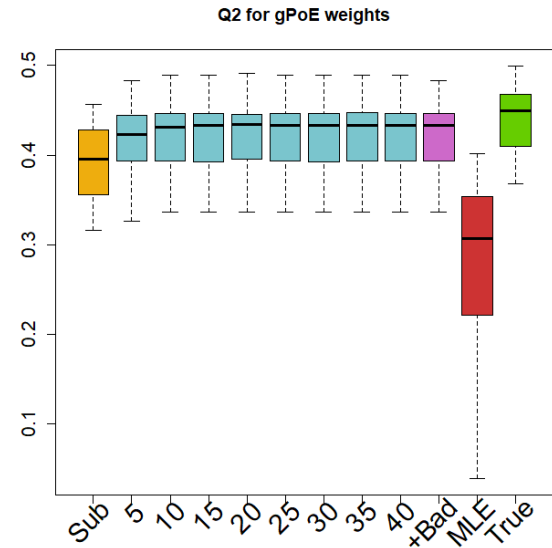
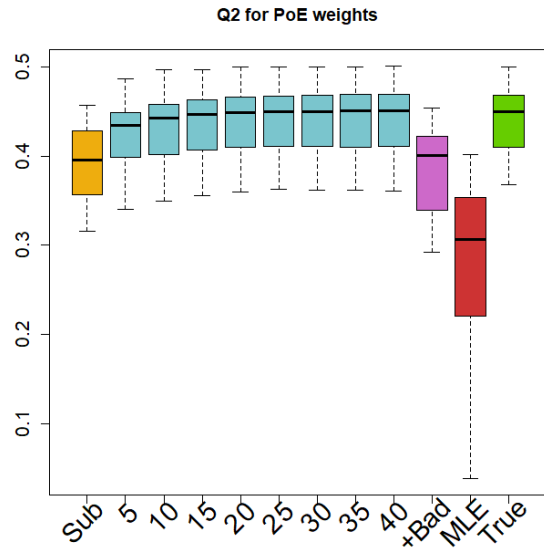
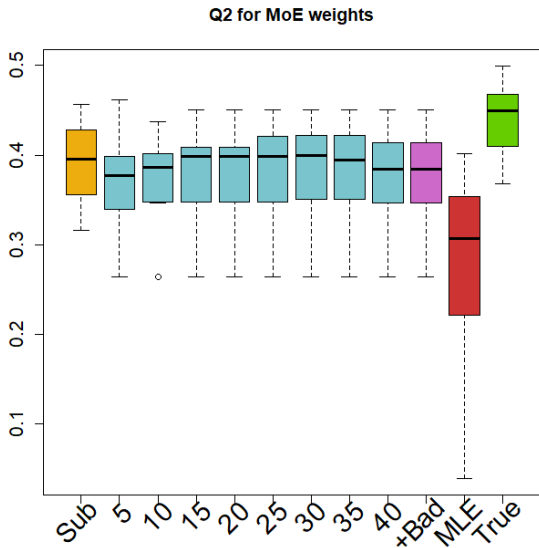
$$f(\theta) \propto \exp(H(R_\theta)).$$

Entropy of a Gaussian correlation in 50D for a uniform design ( $\sigma_X^2 = 1/12$  and  $\kappa_X = 9/5$ ).

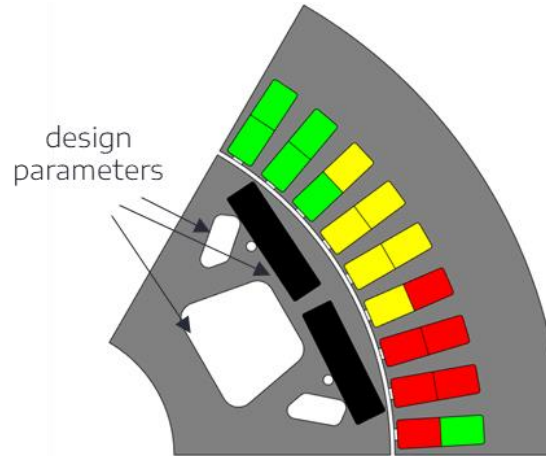
Entropy for a Gaussian correlation



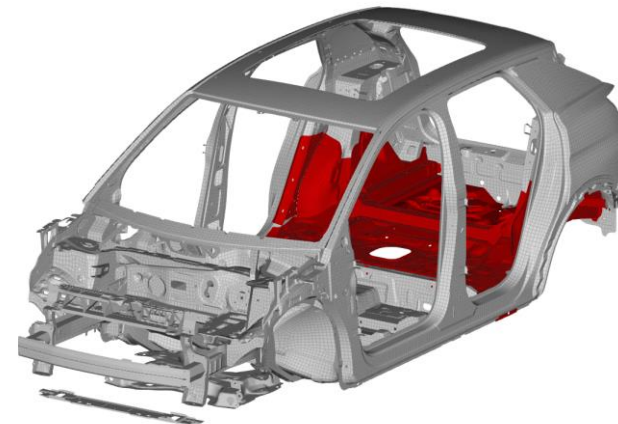
# CHOICE OF THE WEIGHTS



- Résultats sur 2 jeux de données tests
- Étude d'une machine électrique:
  - 37 variables de design,
  - 500 points d'apprentissage,
  - 4500 points test,
  - 2 objectifs et 10 contraintes à modéliser,
  - Résultats moyennés sur 10 runs.

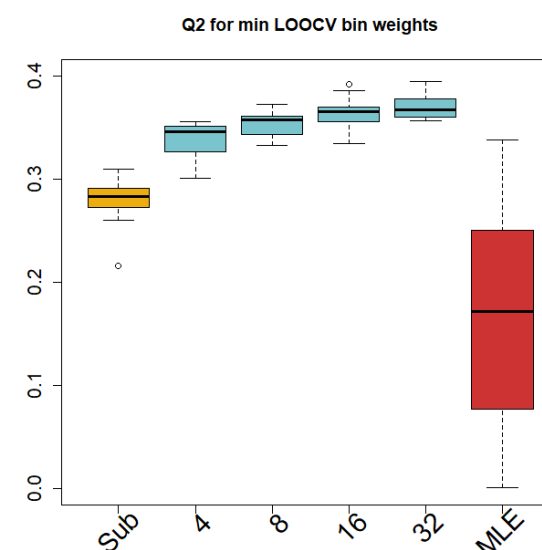
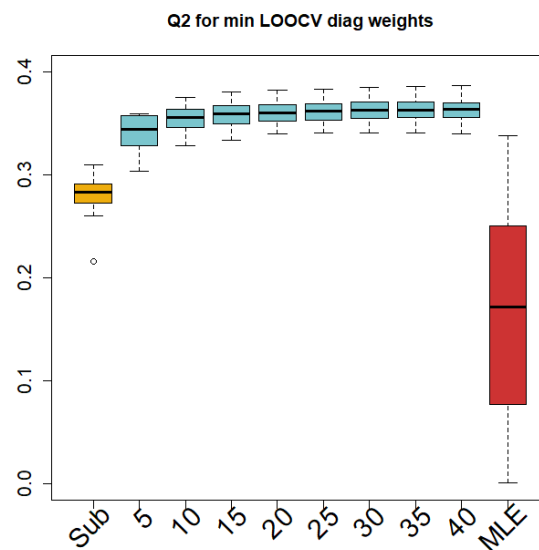
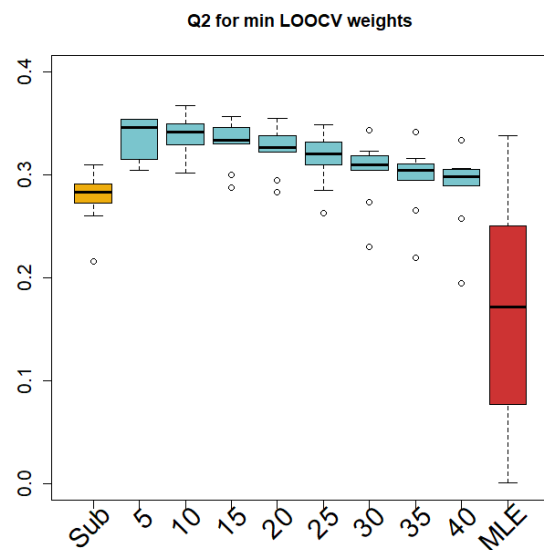
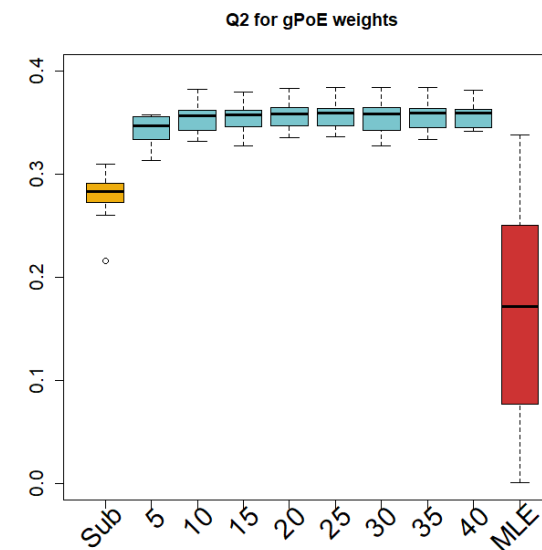
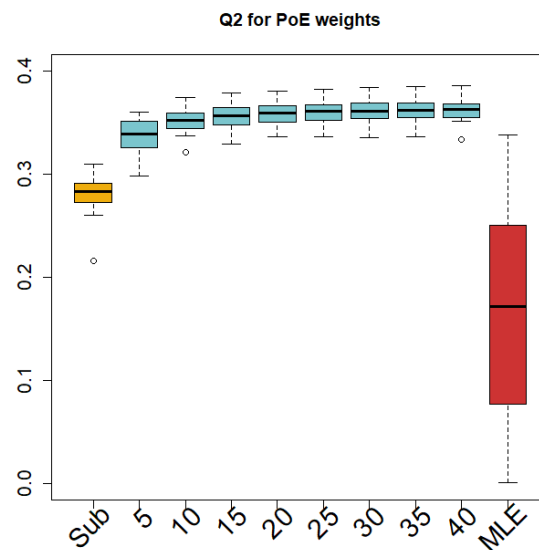
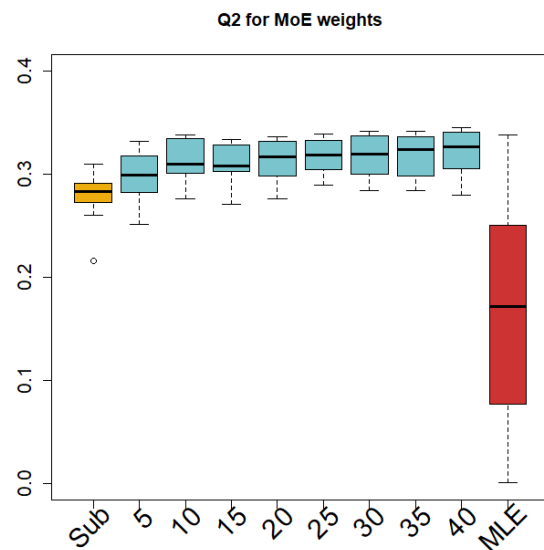


- Étude de la Peugeot 3008 (confort vibratoire et sécurité en crash arrière) :
  - 48 variables de design,
  - 300 points d'apprentissage,
  - 327 points test,
  - 2 objectifs et 413 contraintes (un modèle est construit pour seulement 190 des contraintes).

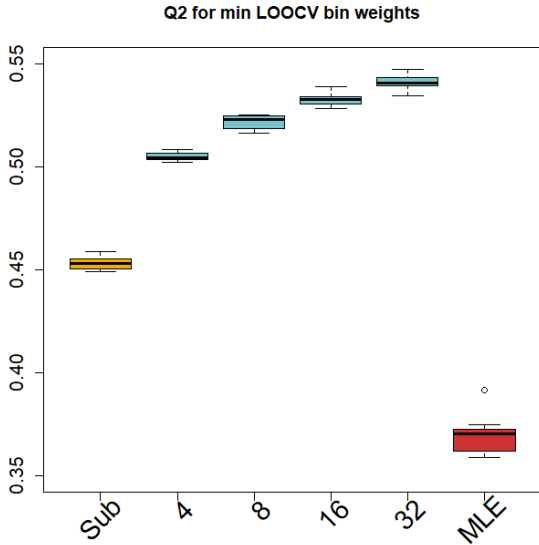
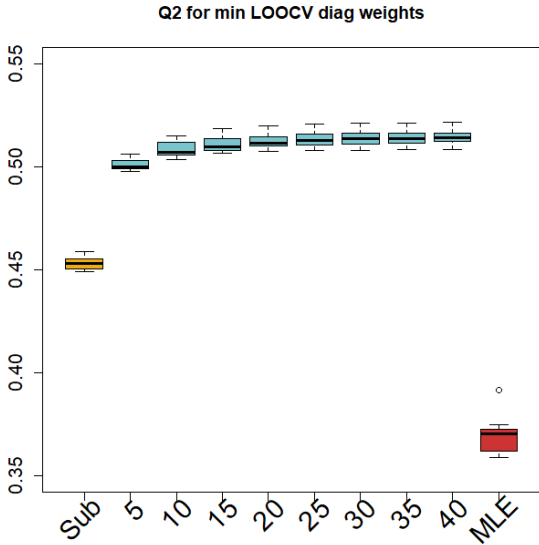
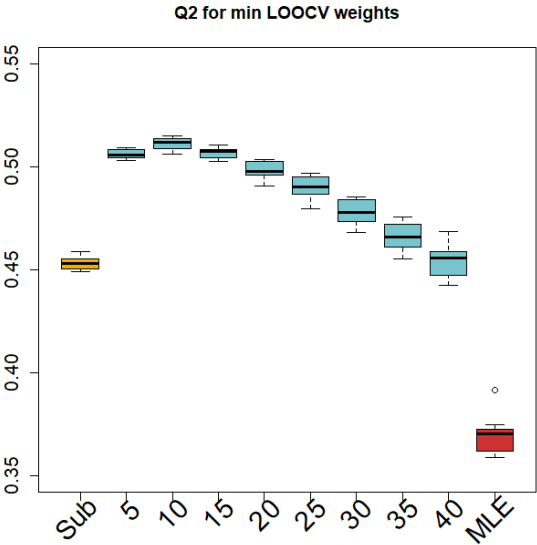
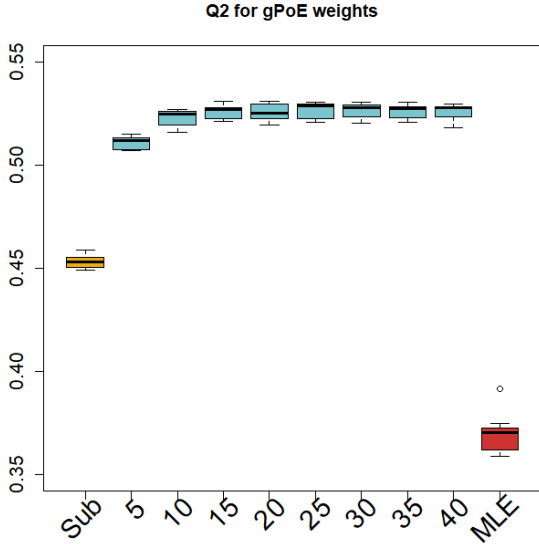
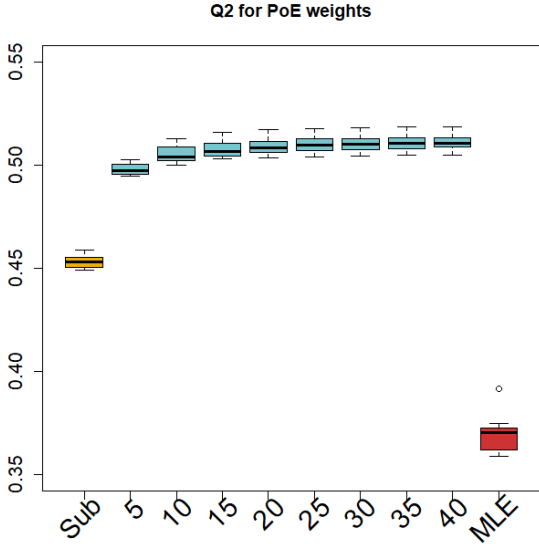
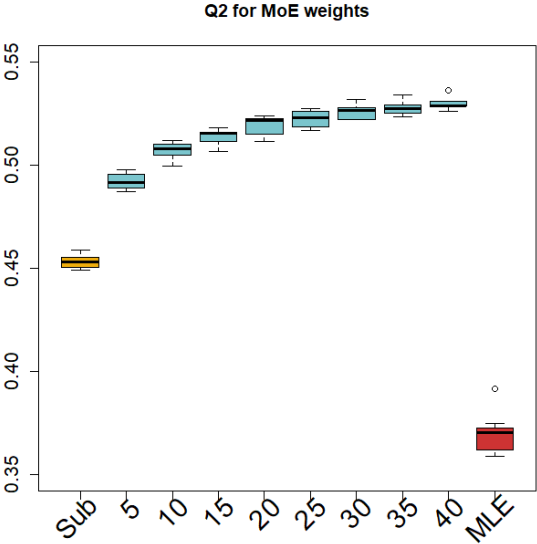




# CHOICE OF THE WEIGHTS



# CHOICE OF THE WEIGHTS





→ We propose a model which is a combination of Kriging models with random length-scales

$$M_{tot}(\mathbf{x}) = \sum_{i=1}^p w_i(\mathbf{x}) M_i(\mathbf{x}),$$

with  $M_i(\mathbf{x}) = \mu_i + k_{\theta_i}(\mathbf{x}, \mathbf{X}) K_{\theta_i}^{-1}(\mathbf{Y} - \mu_i)$  Kriging model with fixed length-scale vector  $\theta_i$ .

### 3) Variance of the combination

## VARIANCE OF THE COMBINATION

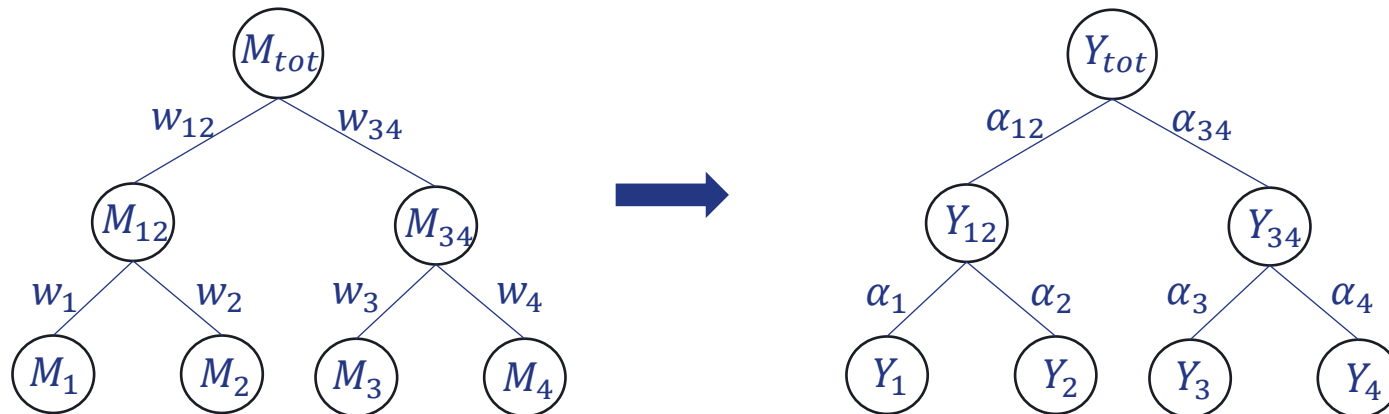
- To obtain the variance of the combination, we add the hypothesis that **the underlying Gaussian Process  $Y_{tot}$  is a combination (with different weights) of independent Gaussian Processes:**

$$Y_{tot} = \sigma_{tot}^2 \sum_{i=1}^p \alpha_i Y_i, \quad \text{with } Y_i \sim \mathcal{GP}(\mu_i, r_{\theta_i}(\dots)), \quad \sum_{i=1}^p \alpha_i = 1, \quad \text{and } \sigma_{tot}^2 \text{ the variance of the GP.}$$

Thus, the covariance of this GP is:

$$k_{tot}(\dots) = \sigma_{tot}^2 \sum_{i=1}^p \alpha_i^2 r_{\theta_i}(\dots).$$

- To simplify the upcoming expressions, we will also assume that the sub-models (and the associated GPs) are combined following a binary tree structure:





- The weights  $\alpha$  in the combination of GPs are chosen to **minimize the expected mean-square error of the combined model** with respect to  $Y_{tot} = \alpha Y_1 + (1 - \alpha)Y_2$  :

$$\alpha^* = \arg \min_{\alpha} \mathbf{E} \left[ \mathbf{E} \left[ \left( wM_1(\mathbf{x}) + (1 - w)M_2(\mathbf{x}) - \alpha Y_1(\mathbf{x}) + (1 - \alpha)Y_2(\mathbf{x}) \right)^2 \mid Y_1, Y_2 \right] \right].$$

By approximation the global MSE using the LOOCV error, we obtain:

$$\alpha^* = \frac{a_1(w)}{a_1(w) + a_2(w)}, \quad \text{with: } \begin{cases} a_1(w) = w^2 \mathbf{E}(e_{LOOCV}(M_1)|Y_2) + (1 - w^2) \mathbf{E}(e_{LOOCV}(M_2)|Y_2), \\ a_2(w) = (1 - w)^2 \mathbf{E}(e_{LOOCV}(M_2)|Y_1) + (1 - (1 - w)^2) \mathbf{E}(e_{LOOCV}(M_1)|Y_1). \end{cases}$$

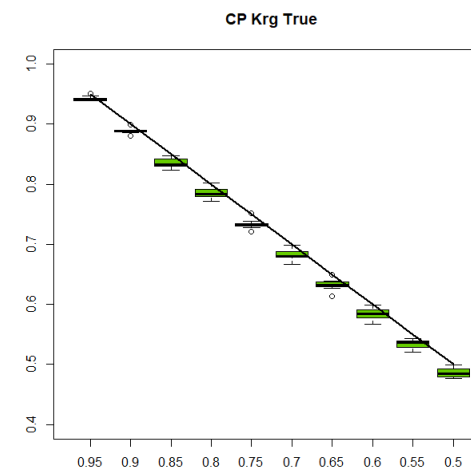
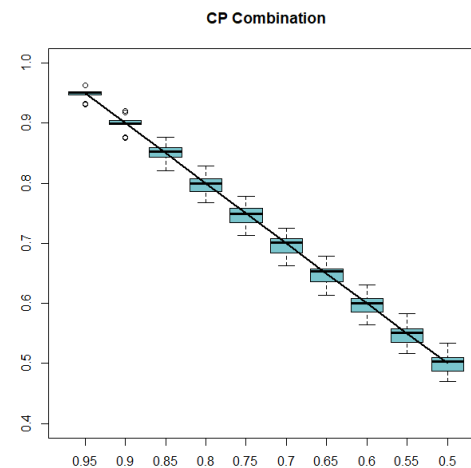
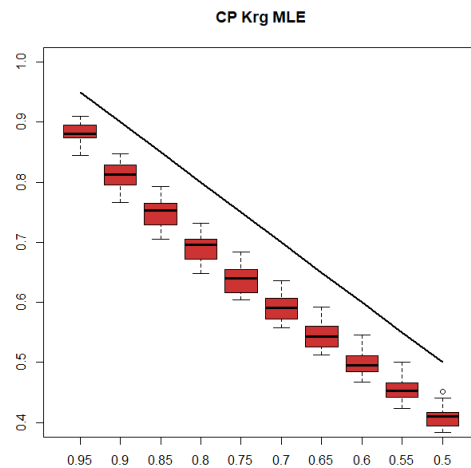
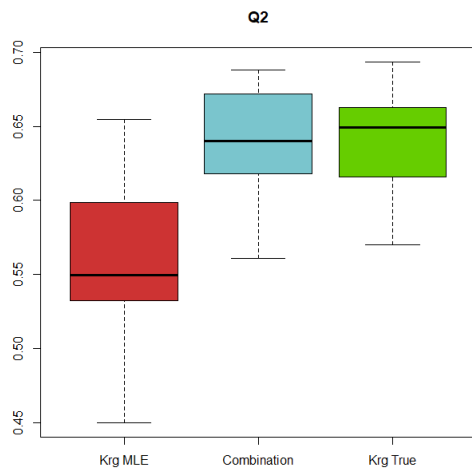
Finally, the variance of the combination is obtained as:

$$\hat{s}^2(\mathbf{x}) = \mathbf{Var}(Y_{tot}(\mathbf{x})|\mathcal{D}) = k_{tot}(\mathbf{x}, \mathbf{x}) - k_{tot}(\mathbf{x}, \mathbf{X})\mathbf{K}_{tot}(\mathbf{X}, \mathbf{X})^{-1}k_{tot}(\mathbf{X}, \mathbf{x}).$$

# RÉSULTATS NUMÉRIQUES – DONNÉES SIMULÉES



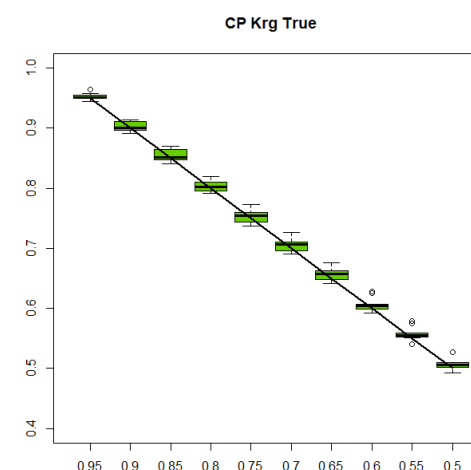
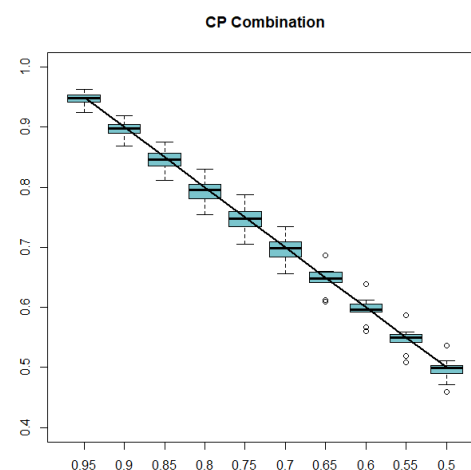
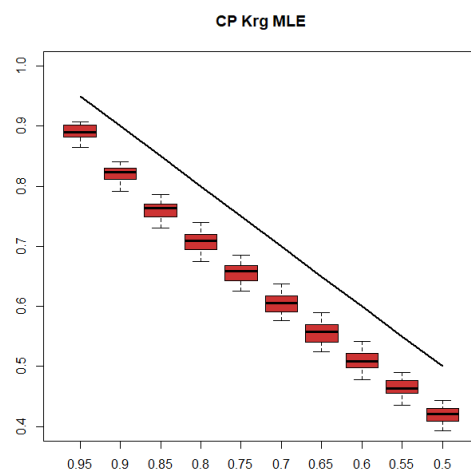
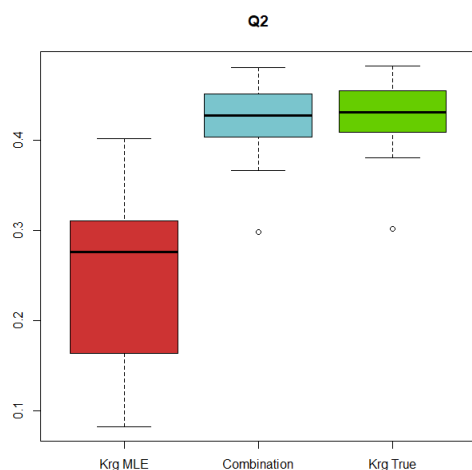
$$\theta_{true} = 3$$



Average computational time:

- Krg MLE: 2,9 mins
- Combination : 0,33 mins

$$\theta_{true} = 2$$

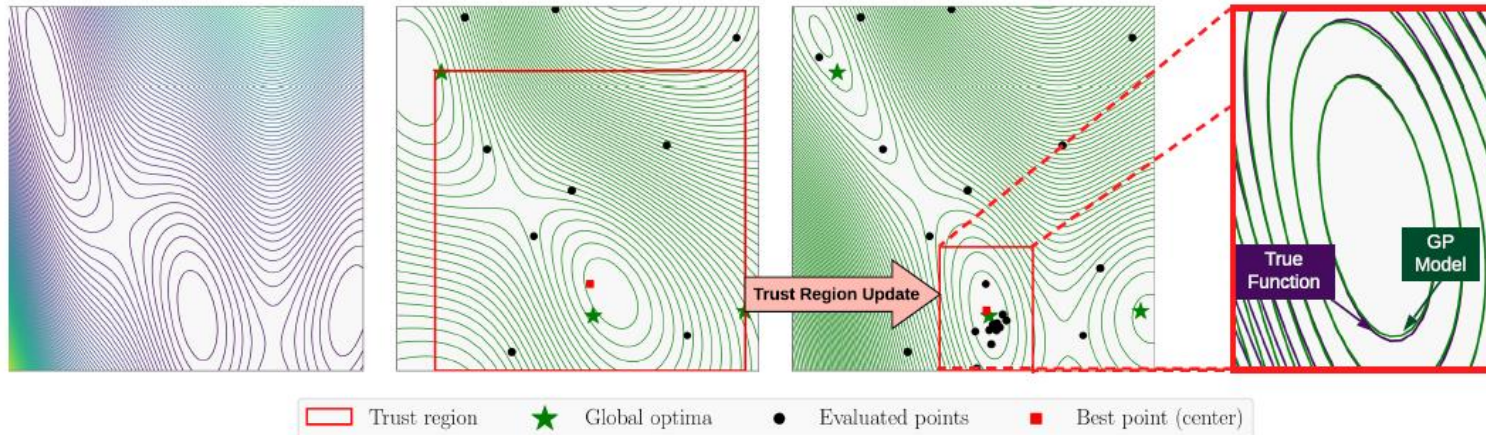


Average computational time:

- Krg MLE: 3,4 mins
- Combination : 0,33 mins

## TRUST REGIONS

- Trust regions are a class of methods to reduce the search space around interesting regions (see for example Eriksson et al., 2019, Diouane et al., 2023).



- A trust region is built around the current best value observed.
- The acquisition function is optimized inside the trust region to select new points.
- After a given number of consecutive failures (no improvement) → The size of the trust region is reduced.
- OR, after a given number of consecutive success (the best value was improved) → The size of the trust region is increased.

## TRUST REGIONS

- In the benchmark, we use the **TREGO** implementation of trust regions (see Diouane et al., 2023).
- In TREGO, we alternate between global iterations of EGO, and local iterations inside a trust region.

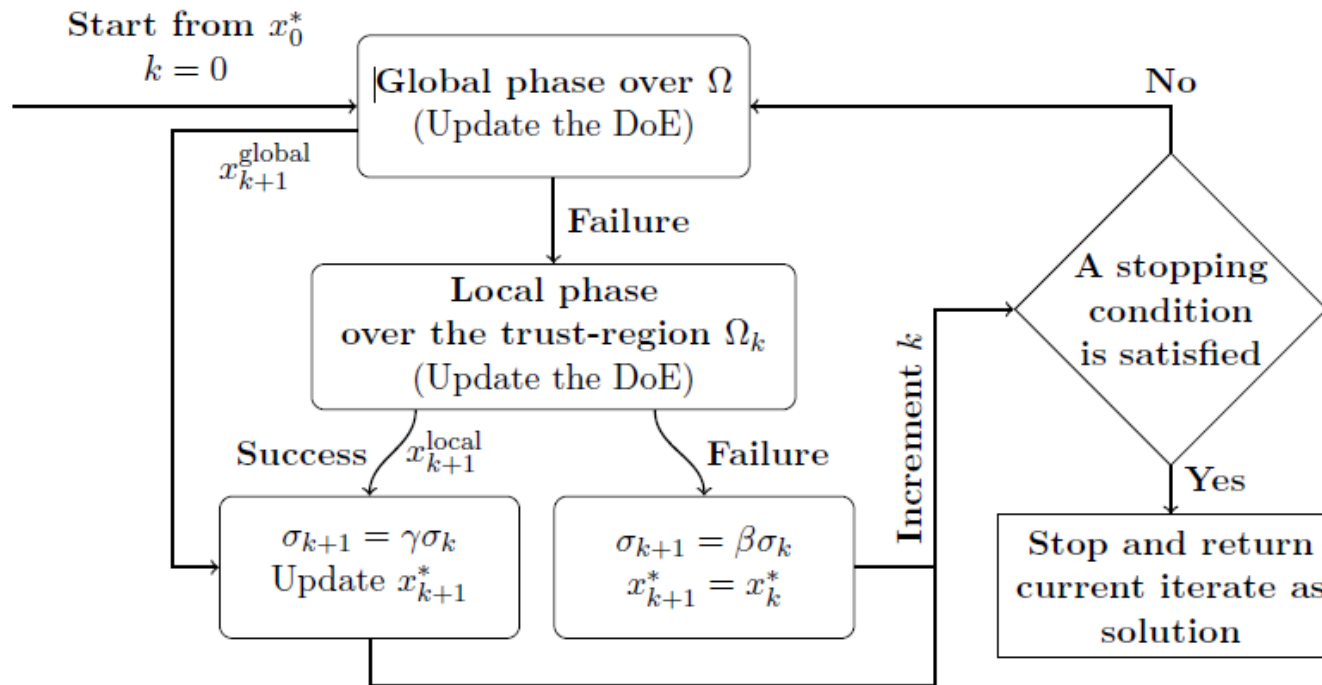


Figure from Diouane et al., 2023 :  
Overview of the TREGO framework.

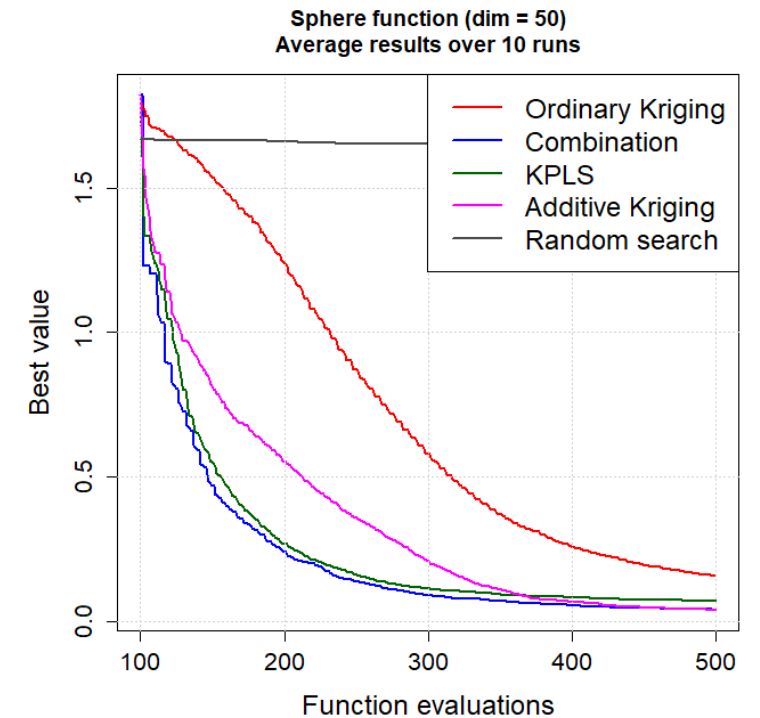
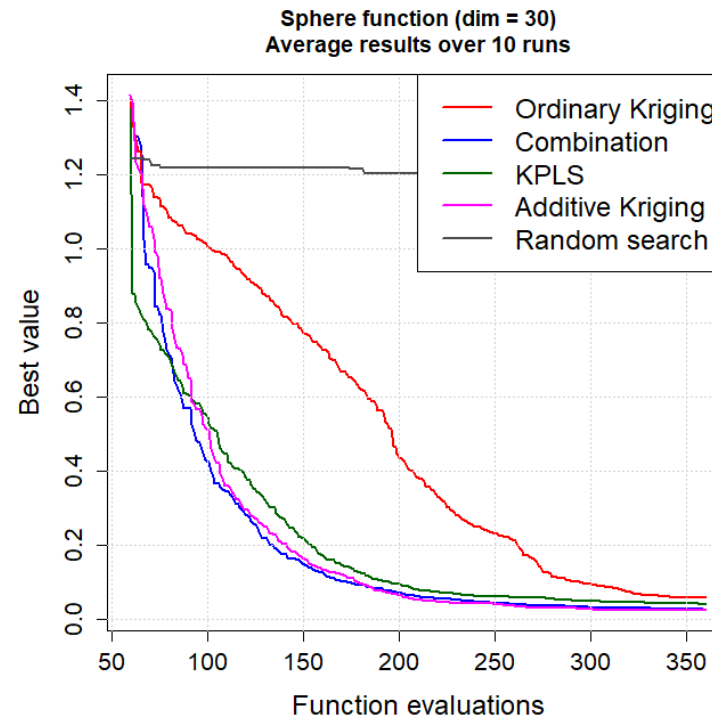
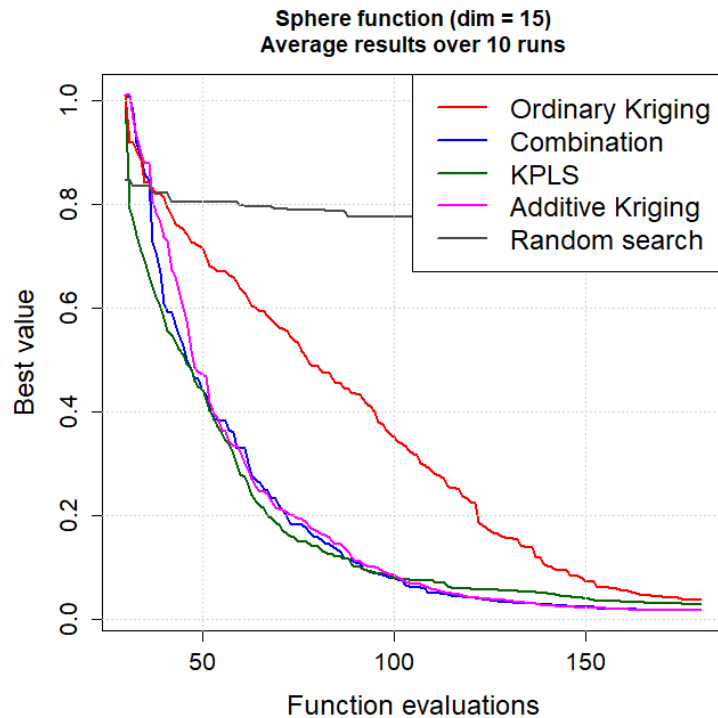
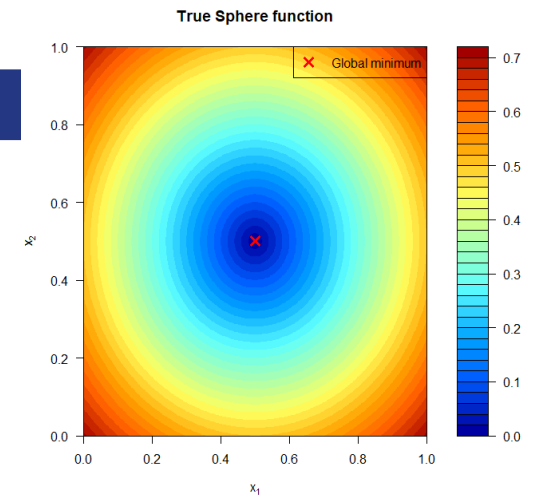


## RÉSULTATS NUMÉRIQUES

- Fonction sphère : 
$$f_{\text{sphère}}(x_1, \dots, x_d) = \sqrt{\sum_{i=1}^d (x_i - 0,5)^2}, \quad 0 \leq x_i \leq 1.$$

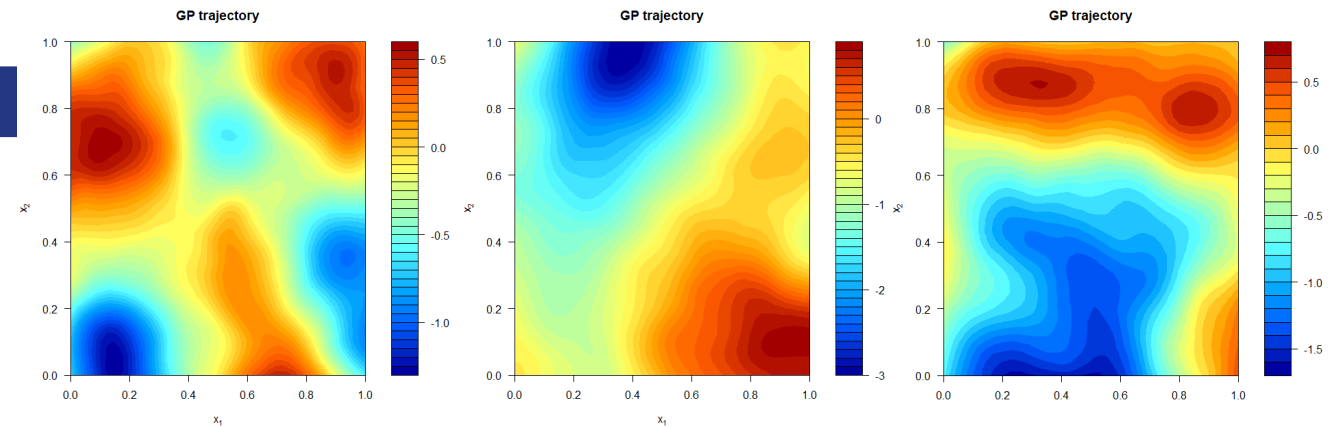
→ Déceptivement difficile à modélisé avec des GPs quand on dispose de peu d'observations.

→ Facile à optimiser (fonction convexe).



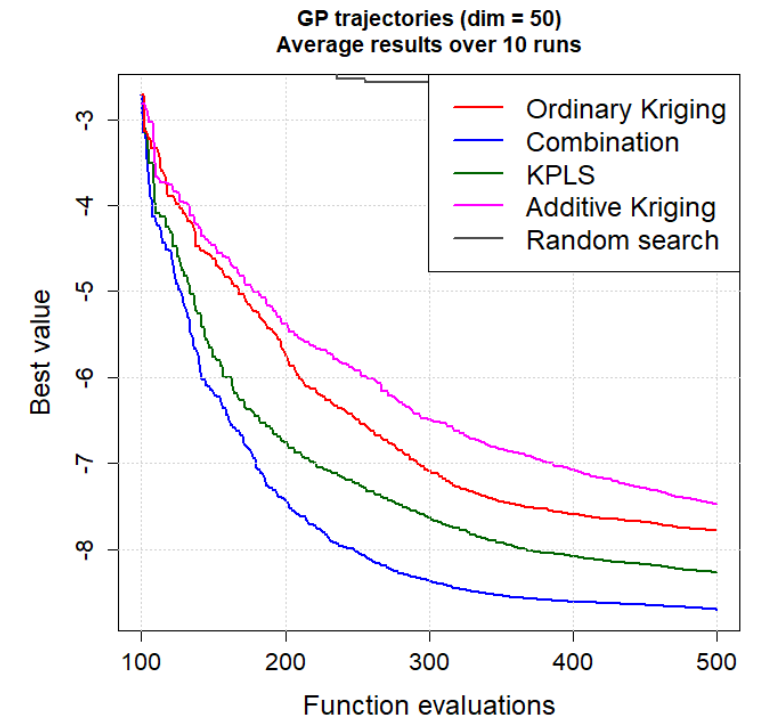
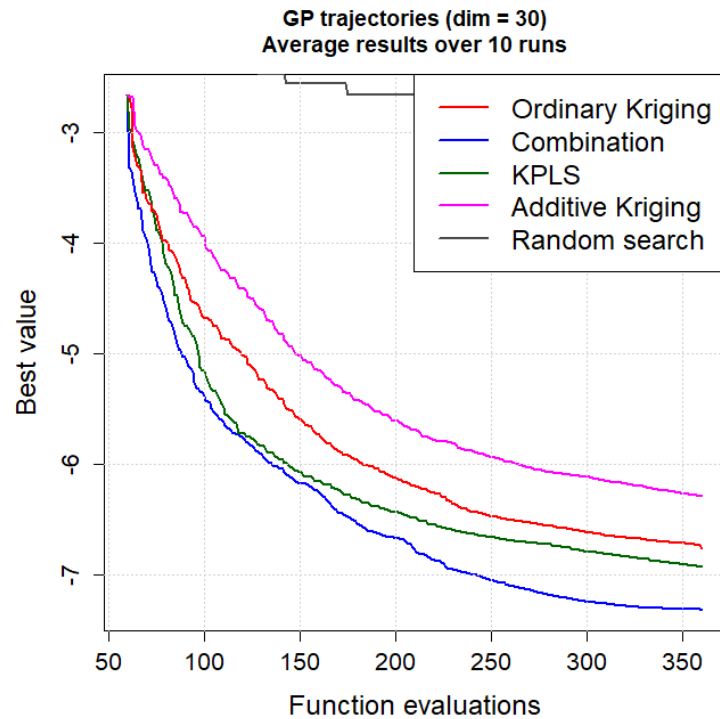
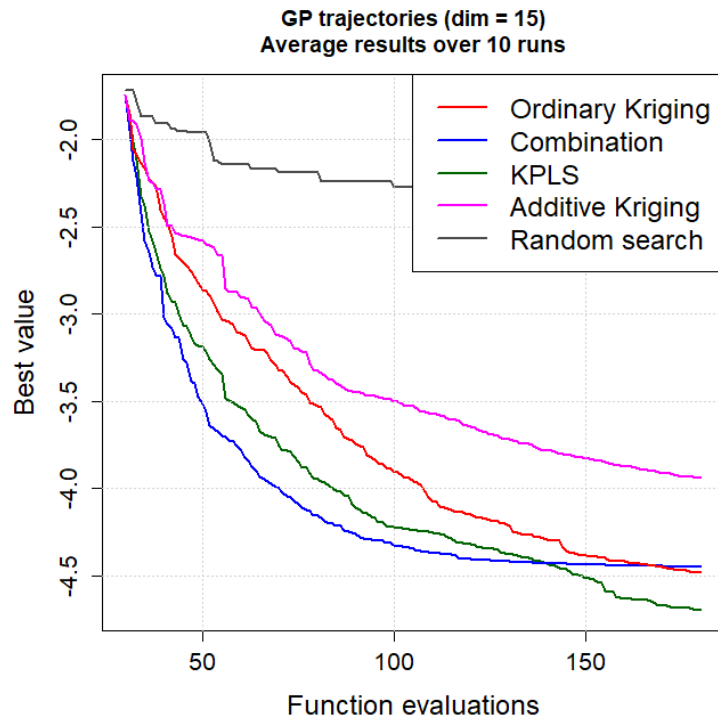
## RÉSULTATS NUMÉRIQUES

- Trajectoires de GPs :  $f_{GP}(\cdot) \sim GP(\mathbf{0}, k_{\theta}(\cdot, \cdot))$ ,  
 $k_{\theta}$  est un noyau Matérn 5/2 isotrope de portée  $\theta = \sqrt{\frac{d}{12}}$



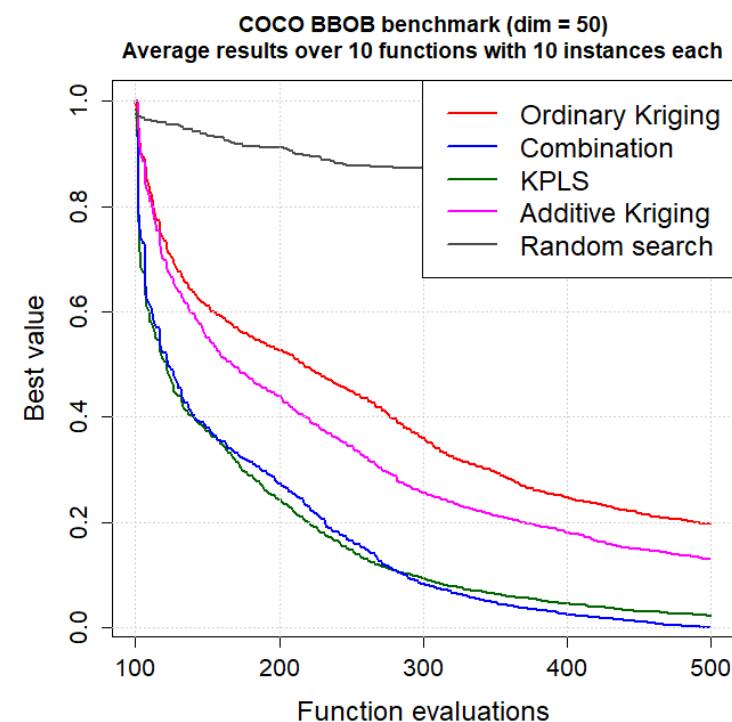
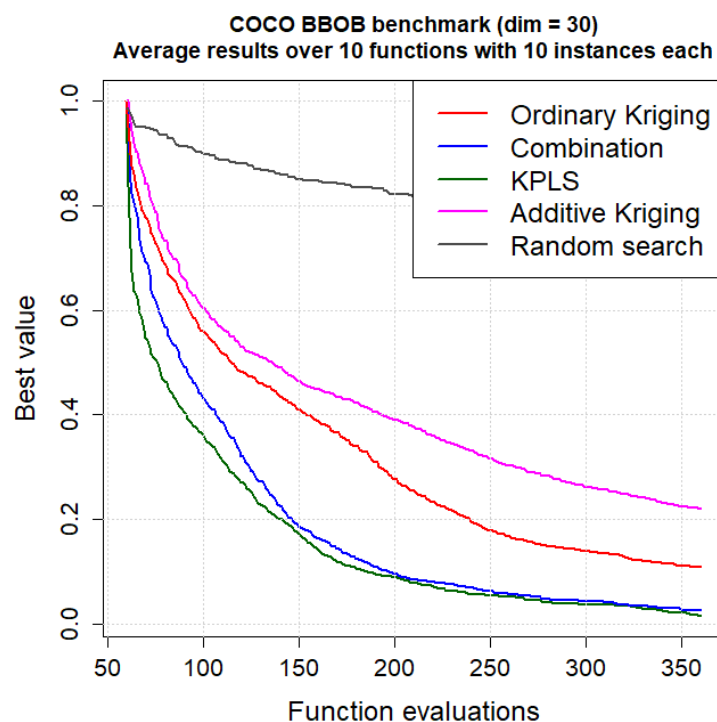
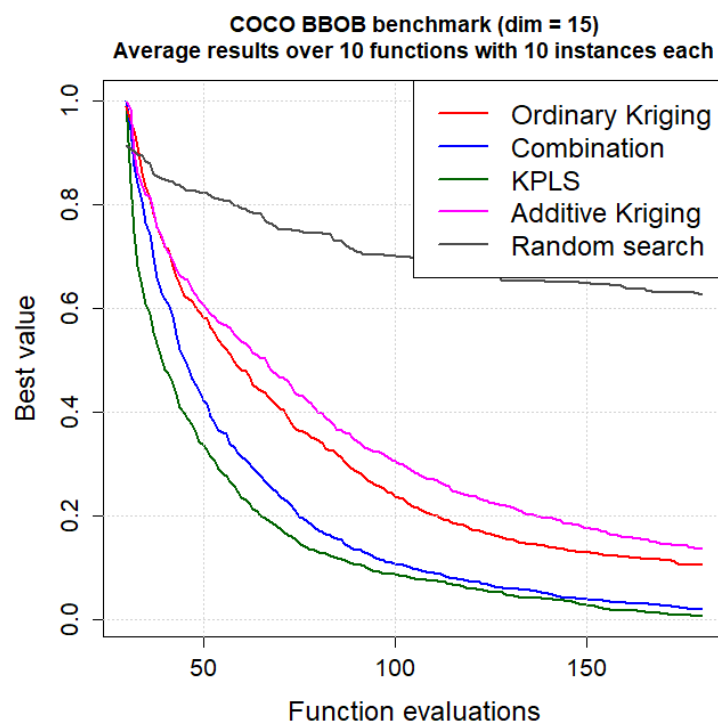
→ Plus difficile à optimiser (multimodale) et plus représentatif des fonctions rencontrées en pratique.

→ Cas où l'hypothèse de Krigeage est vérifiée.



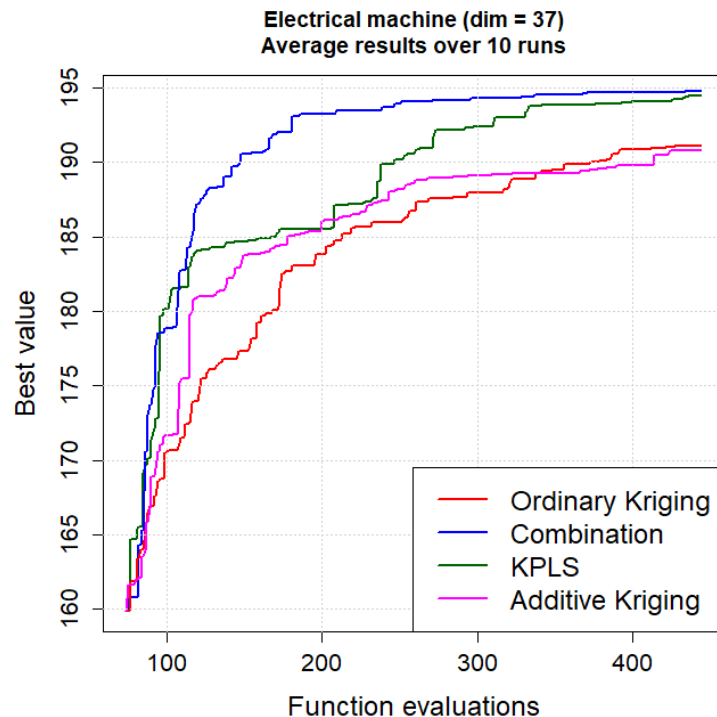
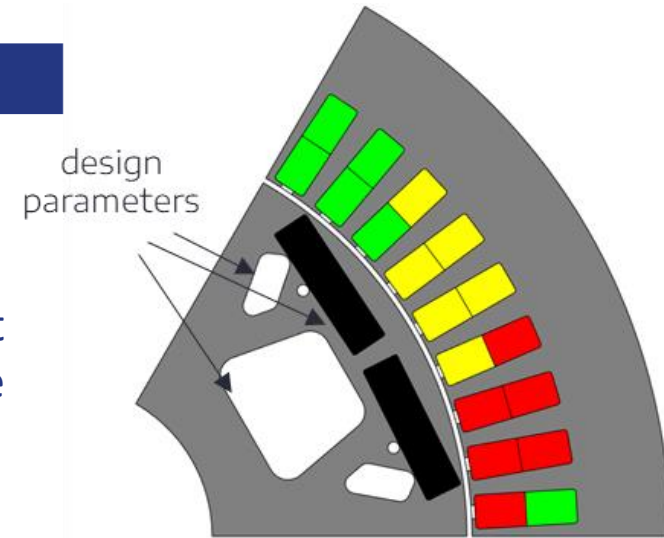
- COCO BBOB benchmark : 9 fonctions multimodales du benchmark (fonctions f15 à f22 et f24)

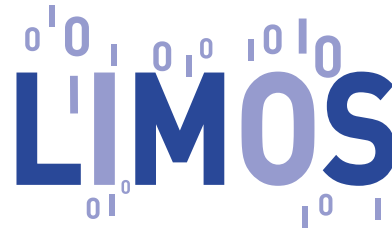
→ Difficile à optimiser (fonctions très multimodales).



## RÉSULTATS NUMÉRIQUES

- Machine électrique (dimension 37)
- Le problème complet comporte 2 objectifs et 10 contraintes. Ici on optimise simplement sur une contrainte (vitesse maximale du véhicule) qui est suffisamment complexe (multimodale) dont de fortes valeurs correspondent globalement à de bonnes machines.





# Thank you for your attention !

Contact :

Tanguy APPRIOU  
[tanguy.appriou@stellantis.com](mailto:tanguy.appriou@stellantis.com)