

# Benchmarking probabilistic spatial machine learning models with complex sample distributions

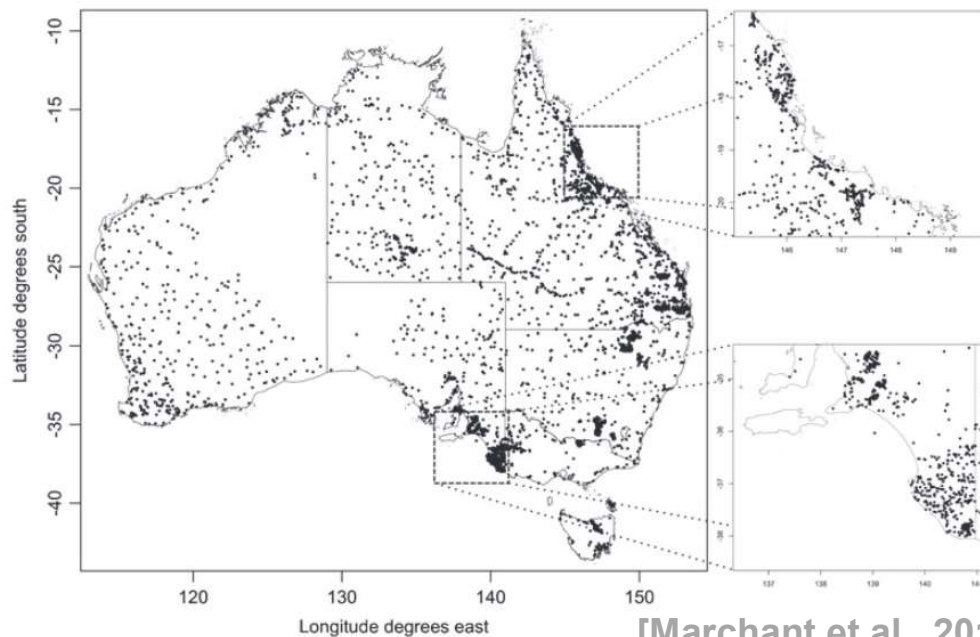
Jeremy Rohmer ([j.rohmer@brgm.fr](mailto:j.rohmer@brgm.fr)), Julie Billy, Vivien Baudouin

4 September 2025



## Examples of complex sample distributions

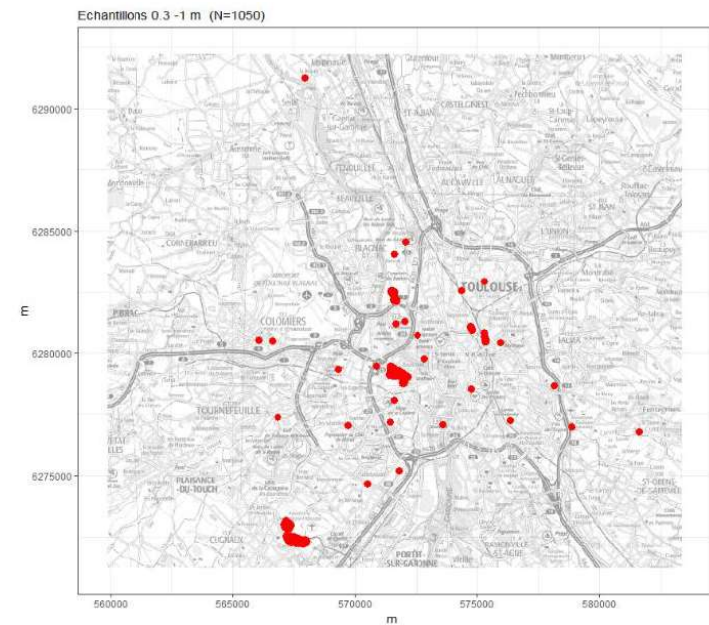
- ❑ **Uneven** spatial distribution with **clusters and sparse samples** in some regions
- ❑ Also owing to **nonstationarities / anisotropies** of the data generating process



[Marchant et al., 2013]

Topsoil samples in Australia

> 2

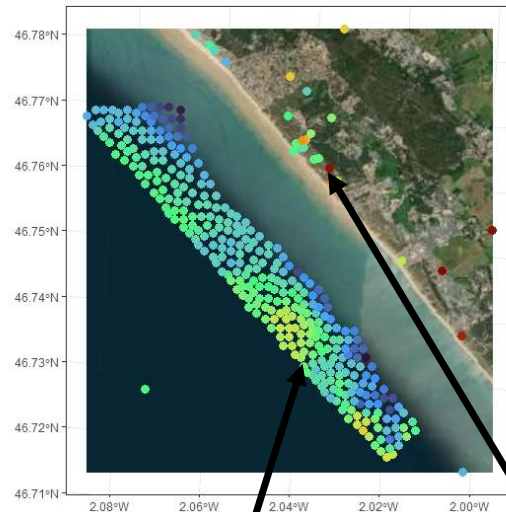


[Belbeze et al., 2019]

Pollutant (Total Petroleum Hydrocarbon) in Toulouse city

> 2

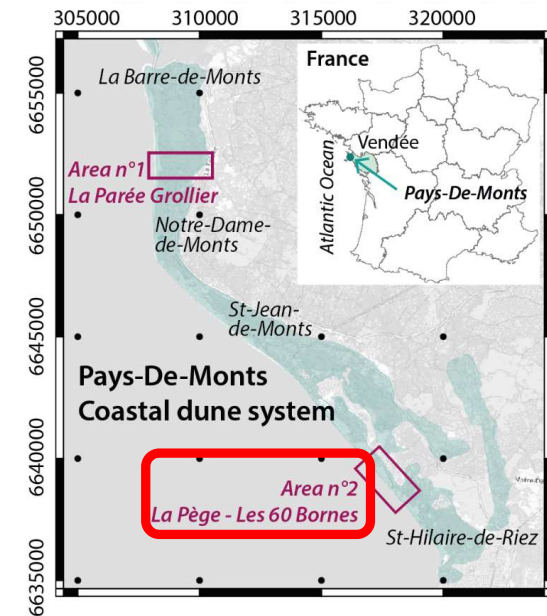
# Examples of complex sample distributions



Highly  
clustered data  
(src: GEOPAL)

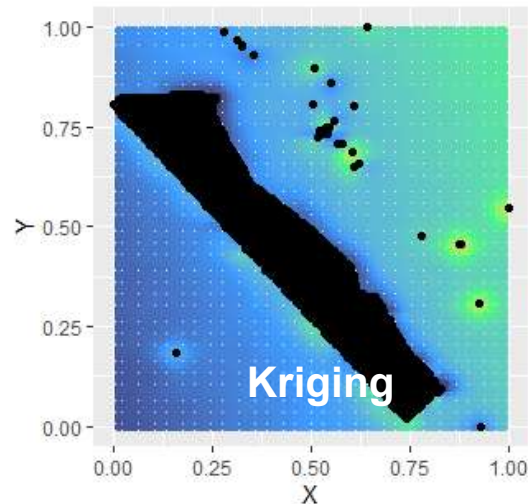


Sparse data  
(src: BSS)

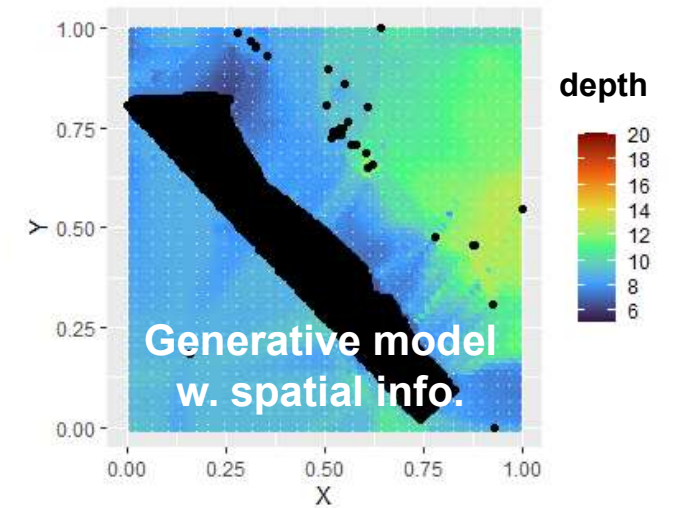
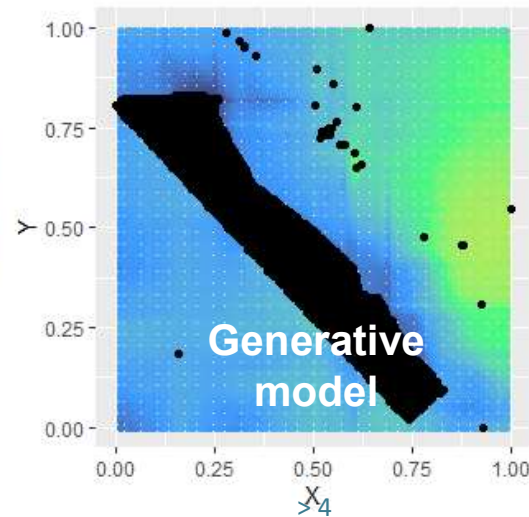
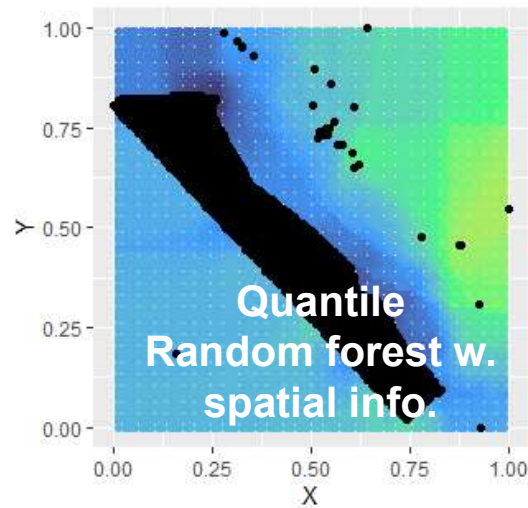
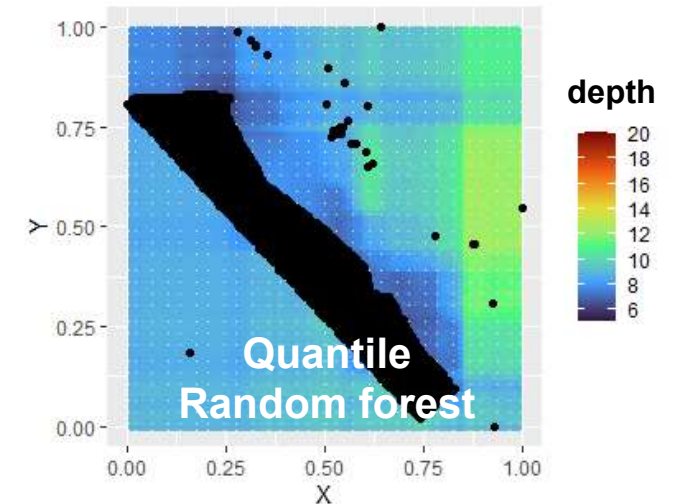


**Interpolation of  
substratum  
topography in the  
dune systems of  
Pays de la Loire**

# Examples of complex sample distributions

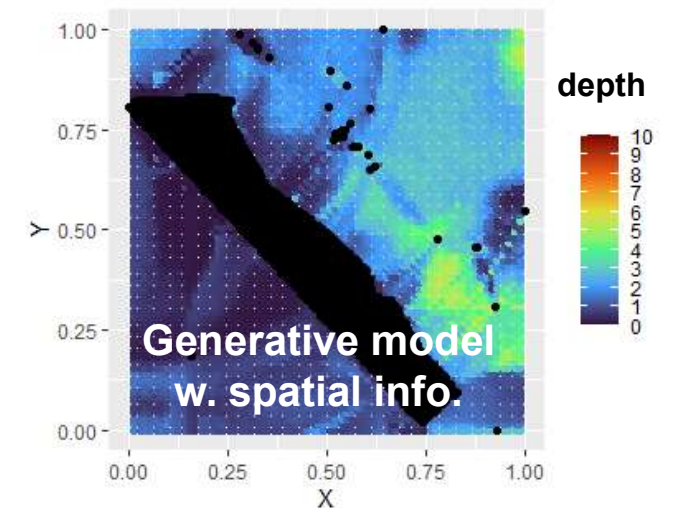
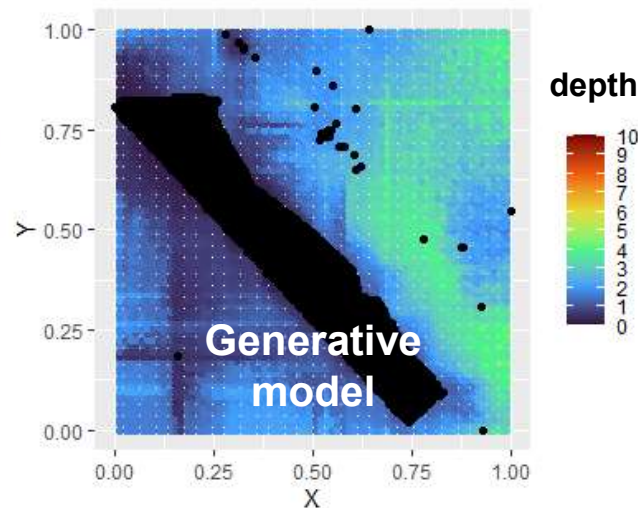
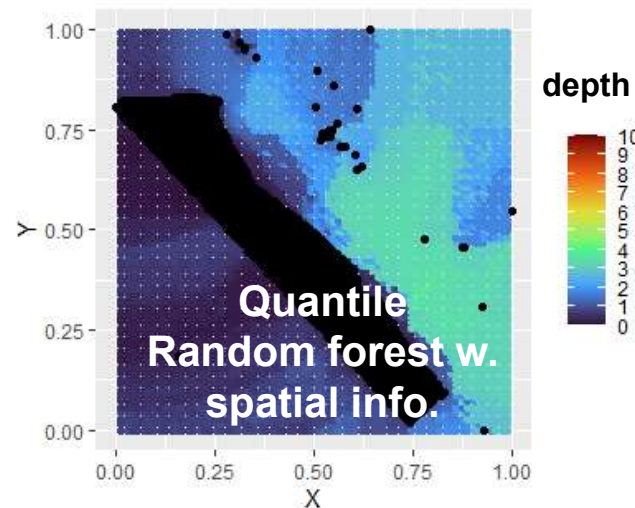
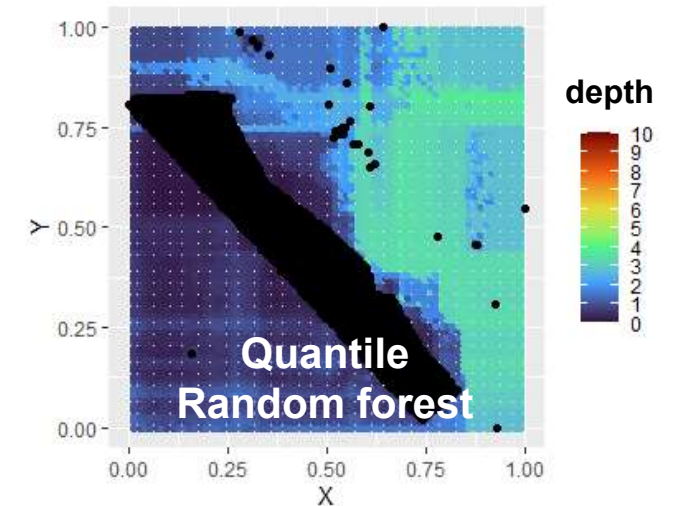
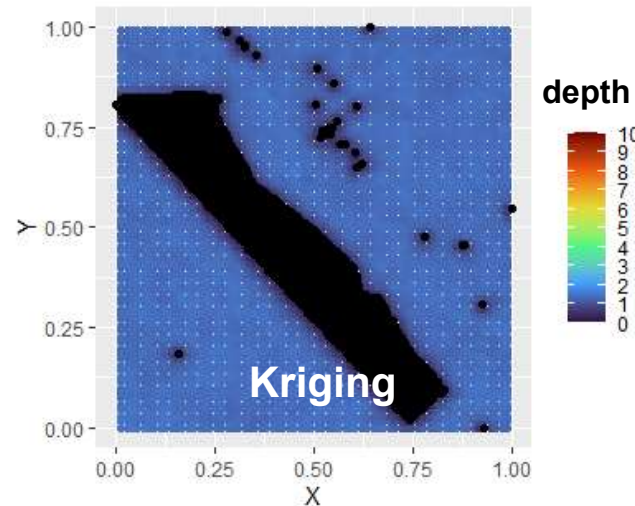


Mean prediction



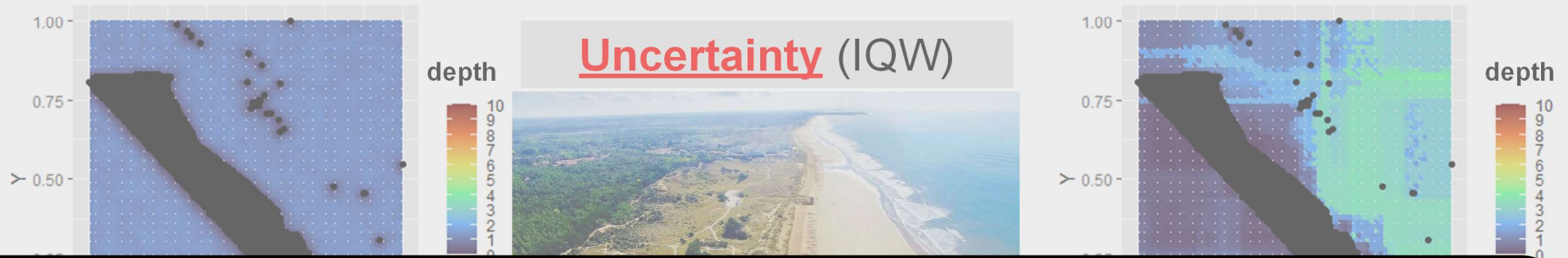


# Examples of complex sample distributions



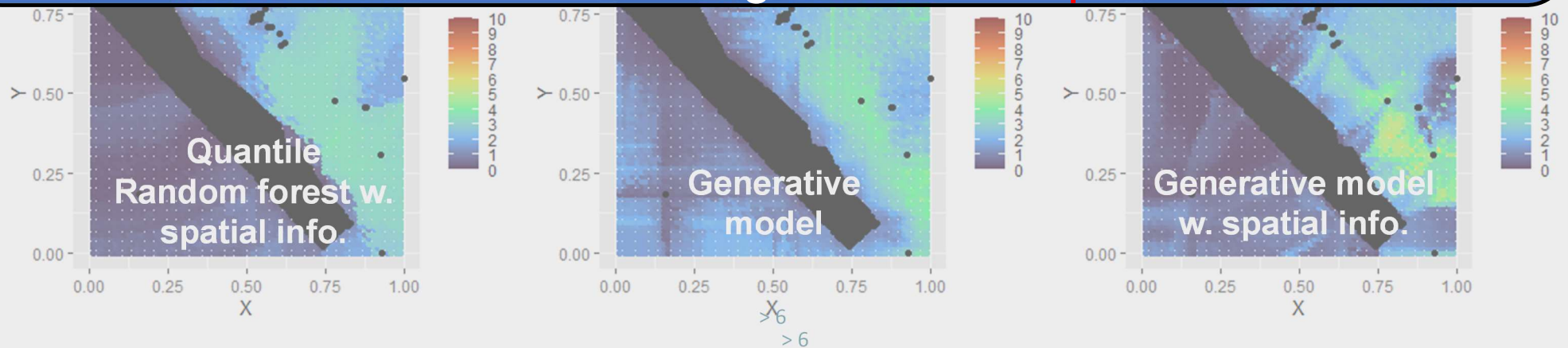


## Motivating real cases



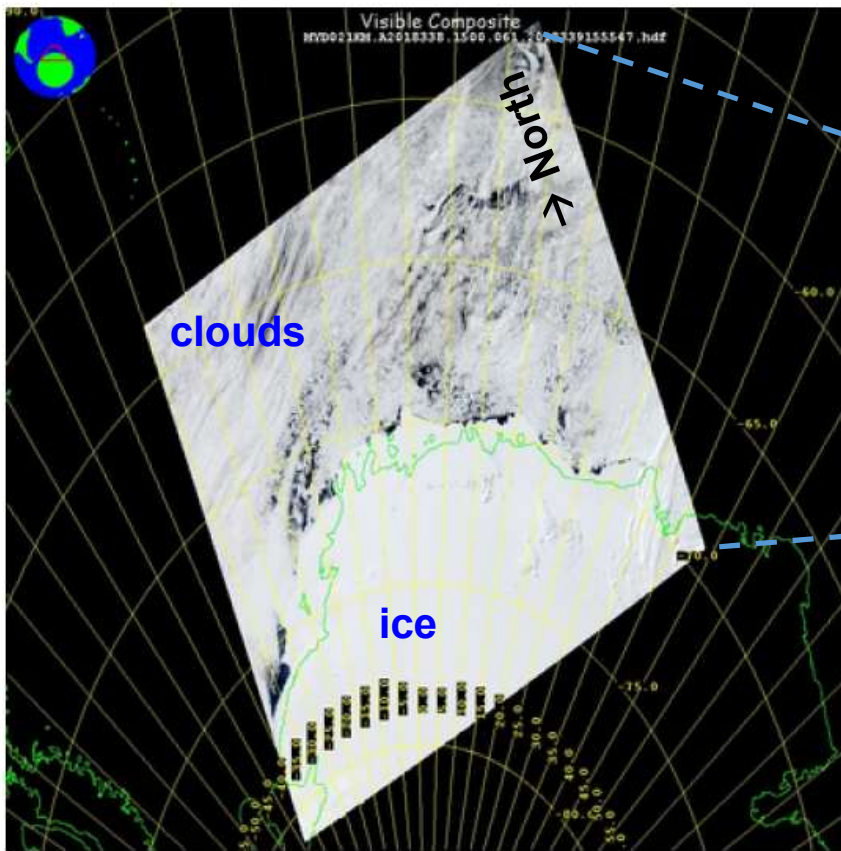
= **Motivation** for a benchmark of **probabilistic ML spatial** models

1. What is the most optimal model(s) ?
2. How to assess the **reliability** of prediction uncertainty?
3. What is the influence of having **clustered / sparse data**?

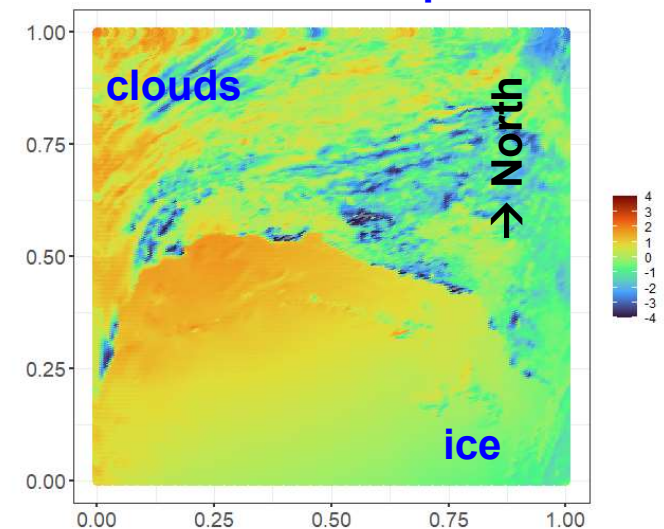
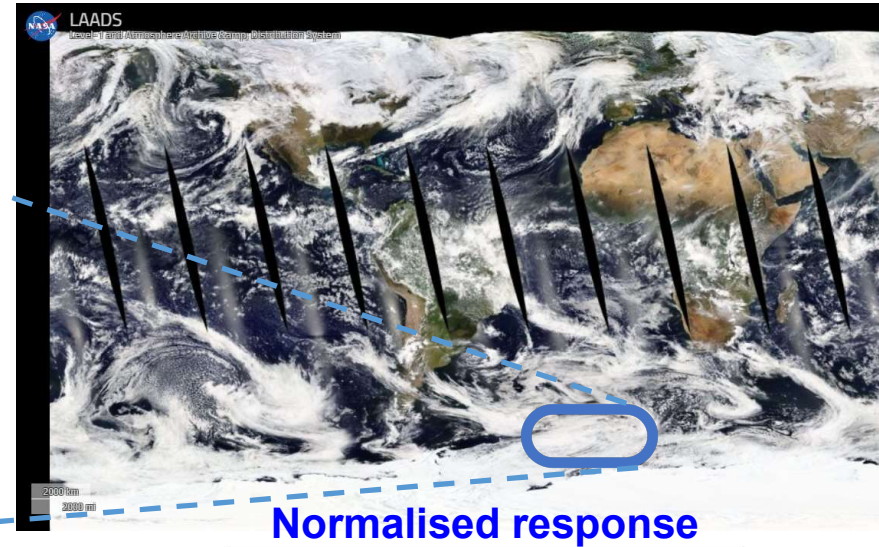




## Benchmark real case with ground truth

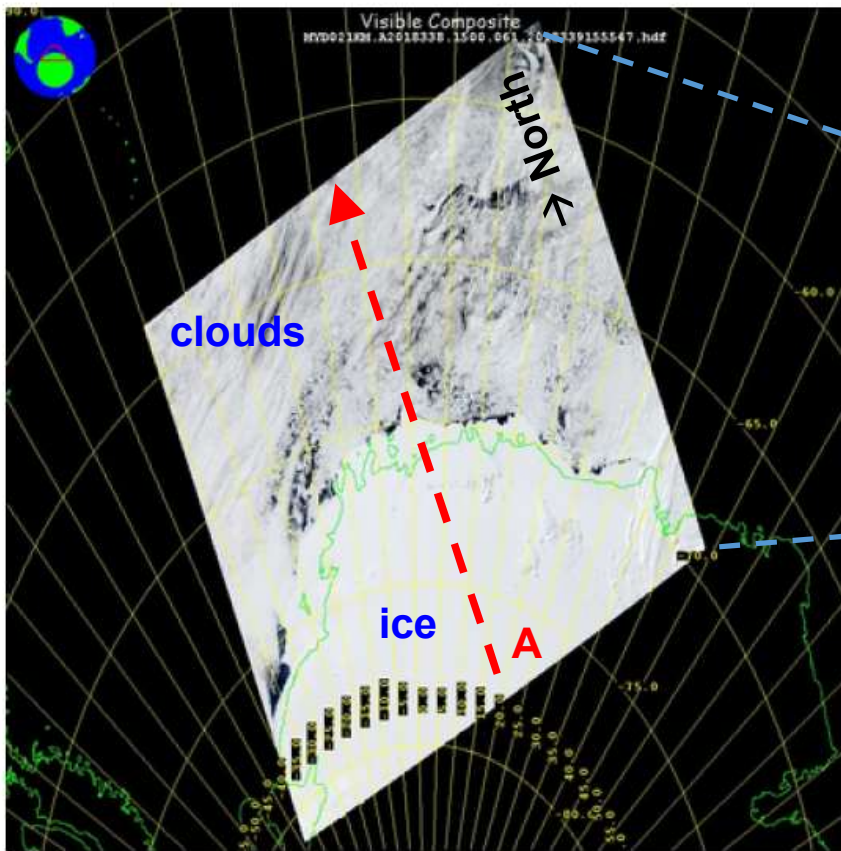


L1B radiances (0:459  $\mu\text{m}$  to 0:479  $\mu\text{m}$  band)  
from the MODIS instrument - Aqua satellite  
(04 December 2018 15:00 UTC)  
extracted from Zammit-Mangion et al. (2022) based on  
<https://ladsweb.modaps.eosdis.nasa.gov>

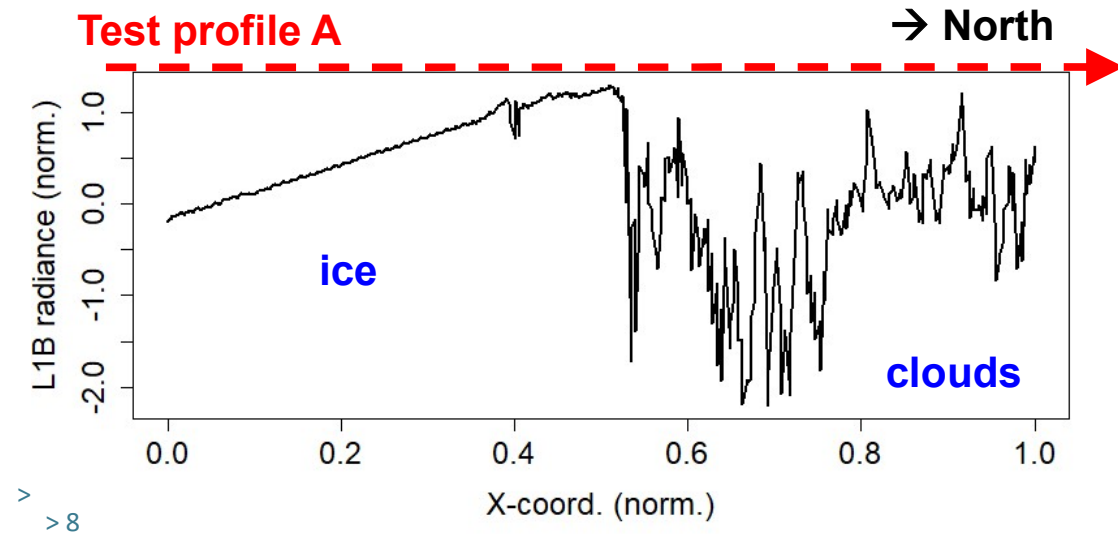
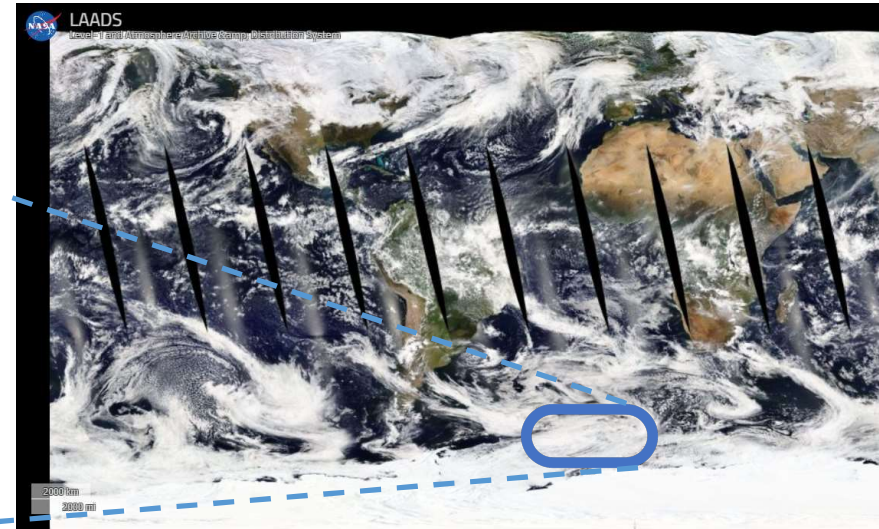




## Benchmark real case with ground truth

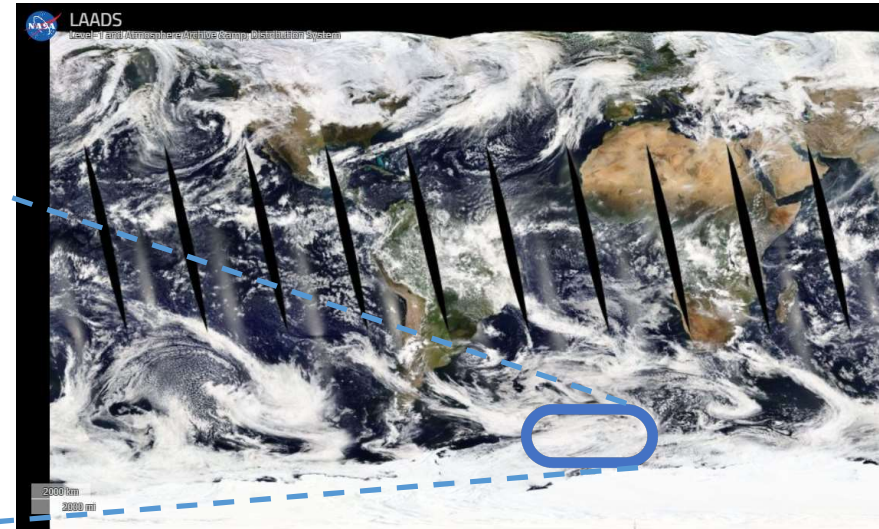
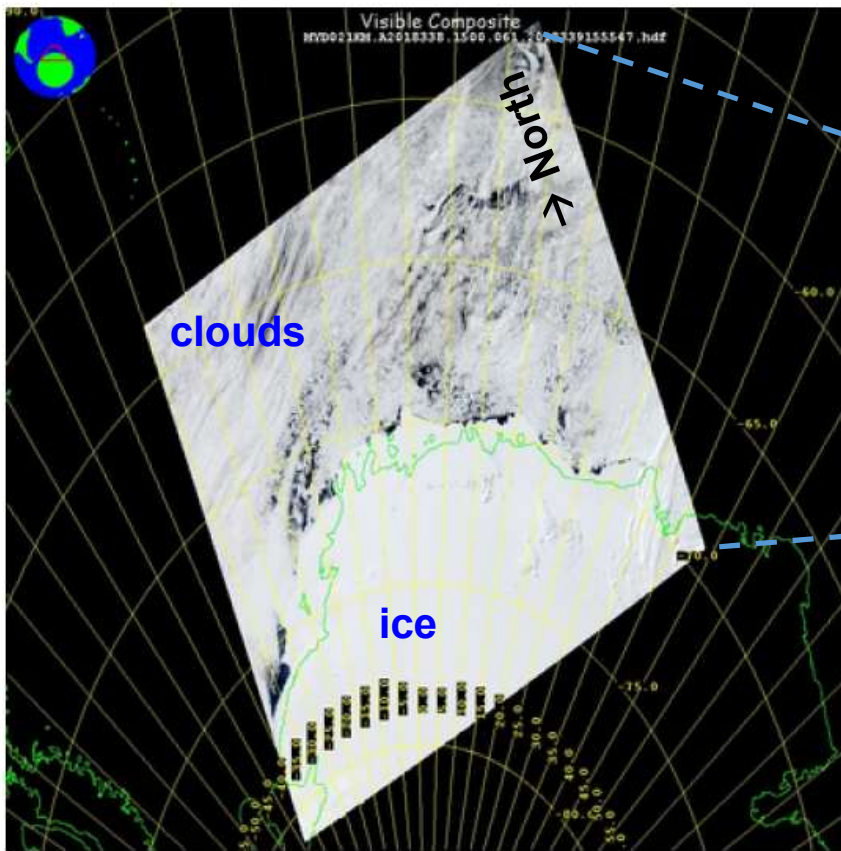


L1B radiances (0:459  $\mu\text{m}$  to 0:479  $\mu\text{m}$  band)  
from the MODIS instrument - Aqua satellite  
(04 December 2018 15:00 UTC)  
extracted from Zammit-Mangion et al. (2022) based on  
<https://ladsweb.modaps.eosdis.nasa.gov>





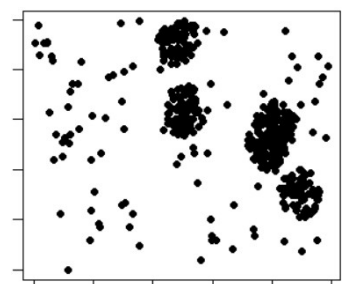
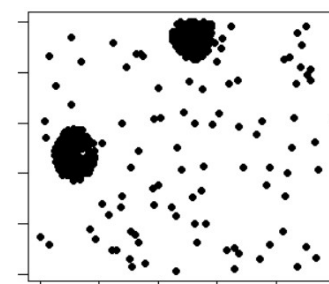
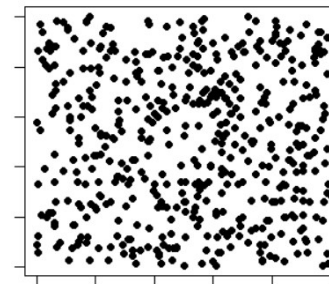
## Benchmark real case with **ground truth**



RANDOM

2 CLUSTERS

4 CLUSTERS



L1B radiances (0:459  $\mu\text{m}$  to 0:479  $\mu\text{m}$  band)  
from the MODIS instrument - Aqua satellite  
(04 December 2018 15:00 UTC)  
extracted from Zammit-Mangion et al. (2022) based on  
<https://ladsweb.modaps.eosdis.nasa.gov>

$N=500$ ,  $N_o=20\%$  of samples outside the clustered regions  
2D Covariates = **spatial coordinates**



## Performance **scores**

Define the test set  $\mathbf{T} = (\mathbf{X}_i, y_i)_{i=1, \dots, n}$  where the response  $Y$  is related to spatial coordinates  $\mathbf{X}$

□ Measure of **accuracy**: coefficient of determination

$$Q^2 = 1 - \frac{\sum_{i \in \mathbf{T}} (y_i - \hat{\mu}_i)^2}{\sum_{i \in \mathbf{T}} (y_i - \bar{y})^2} \quad \text{where } \hat{\mu} \text{ is the ML conditional mean}$$

**Compared to**  **1.0**



## Performance **scores**

Define the test set  $\mathbf{T} = (\mathbf{X}_i, y_i)_{i=1, \dots, n}$  where the response  $Y$  is related to spatial coordinates  $\mathbf{X}$

□ Measure of **accuracy**: coefficient of determination

$$Q^2 = 1 - \frac{\sum_{i \in \mathbf{T}} (y_i - \hat{\mu}_i)^2}{\sum_{i \in \mathbf{T}} (y_i - \bar{y})^2} \quad \text{where } \hat{\mu} \text{ is the ML conditional mean} \quad \xrightarrow{\text{Compared to}} \quad 1.0$$

□ Measure of **'statistical' accuracy** (calibration): coverage score for **prediction interval**  $PI^\alpha = [\hat{Q}^{\alpha/2}; \hat{Q}^{1-\alpha/2}]$

$$Cov = \frac{1}{|\mathbf{T}|} \sum_{i \in \mathbf{T}} \mathbf{1}(y_i \in PI^\alpha) \quad \text{where } \hat{Q} \text{ is the ML conditional quantile} \quad \xrightarrow{\text{Compared to}} \quad 1 - \alpha$$



## Performance **scores**

Define the test set  $\mathbf{T} = (\mathbf{X}_i, y_i)_{i=1, \dots, n}$  where the response  $Y$  is related to spatial coordinates  $\mathbf{X}$

□ Measure of **accuracy**: coefficient of determination

$$Q^2 = 1 - \frac{\sum_{i \in T} (y_i - \hat{\mu}_i)^2}{\sum_{i \in T} (y_i - \bar{y})^2} \quad \text{where } \hat{\mu} \text{ is the ML conditional mean} \quad \xrightarrow{\text{Compared to}} \quad 1.0$$

□ Measure of **'statistical' accuracy** (calibration): coverage score for **prediction interval**  $PI^\alpha = [\hat{Q}^{\alpha/2}; \hat{Q}^{1-\alpha/2}]$

$$Cov = \frac{1}{|\mathbf{T}|} \sum_{i \in T} \mathbf{1}(y_i \in PI^\alpha) \quad \text{where } \hat{Q} \text{ is the ML conditional quantile} \quad \xrightarrow{\text{Compared to}} \quad 1 - \alpha$$

□ Measure (weighted) **informativeness** of  $PI^\alpha$ : interval score [Gneiting & Raftery 2007]

$$IS_i^\alpha = \underbrace{(\hat{Q}^{1-\alpha/2} - \hat{Q}^{\alpha/2})}_{\text{sharpness}} + \underbrace{\frac{2}{\alpha}(\hat{Q}^{\alpha/2} - y_i)\mathbf{1}(y_i < \hat{Q}^{\alpha/2})}_{\text{underprediction}} + \underbrace{\frac{2}{\alpha}(y_i - \hat{Q}^{1-\alpha/2})\mathbf{1}(y_i > \hat{Q}^{1-\alpha/2})}_{\text{overprediction}}$$

$$\xrightarrow{\text{Compared to}} \quad 0.0$$

## Class 1 of spatial probabilistic ML models: GP-like

### □ Gaussian process regression ('typical / shallow' GP)

Conditioned on the data points  $(\mathbf{X}_i, Y_i)_{i=1,\dots,n}$  where the response  $Y$  is related to spatial coordinates  $\mathbf{X}$

$$Y(\mathbf{X}^*) \sim \text{Gauss}(\mu^*, C^*)$$

where the conditional  $\mu^*, C^*$  are given by the 'typical' kriging equations from  $\mathbf{X}, Y$  [Rasmussen & Williams 2006]



## Class 1 of spatial probabilistic ML models: GP-like

### □ Gaussian process regression ('typical' / shallow' GP)

Conditioned on the data points  $(X_i, Y_i)_{i=1,\dots,n}$  where the response  $Y$  is related to spatial coordinates  $X$

$$Y(X^*) \sim \text{Gauss}(\mu^*, C^*)$$

where the conditional  $\mu^*, C^*$  are given by the 'typical' kriging equations from  $X, Y$  [Rasmussen & Williams 2006]

### □ Deep Gaussian process (DGP):

Successive warping (special case of nested GPs) to handle nonstationarities [Wikle & Zammit-Mangion 2022]

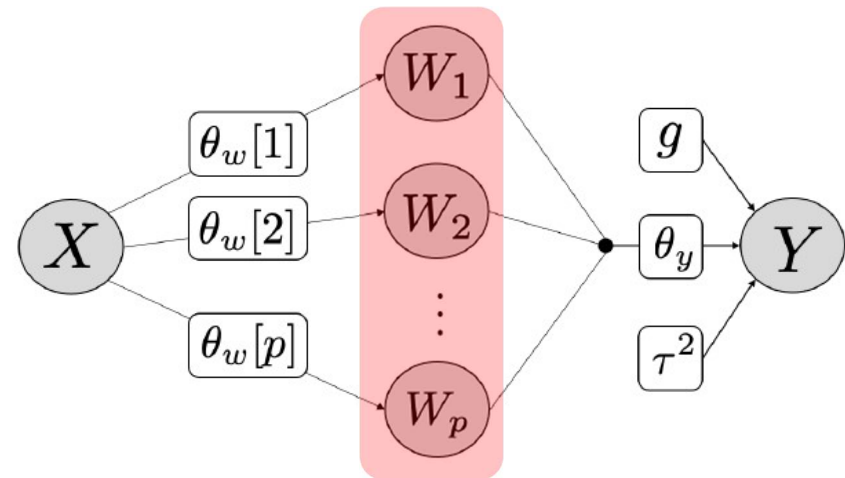
$$Y(X^*) | \mathbf{W} \sim \text{Gauss}(0, C(\mathbf{W}))$$

$$\mathbf{W}_k \sim \text{Ind} \text{Gauss}(0, C(X)) \quad \forall k = 1, \dots, p$$

#### Assumptions

- Latent GP  $\mathbf{W}$  unit scale, noise free
- Conditional independence among nodes of  $\mathbf{W}$
- Isotropic lengths  $\theta$

Full Bayesian inference using MCMC scheme combined with Elliptical slice sampling for  $\mathbf{W}$  [Sauer et al., 2022]

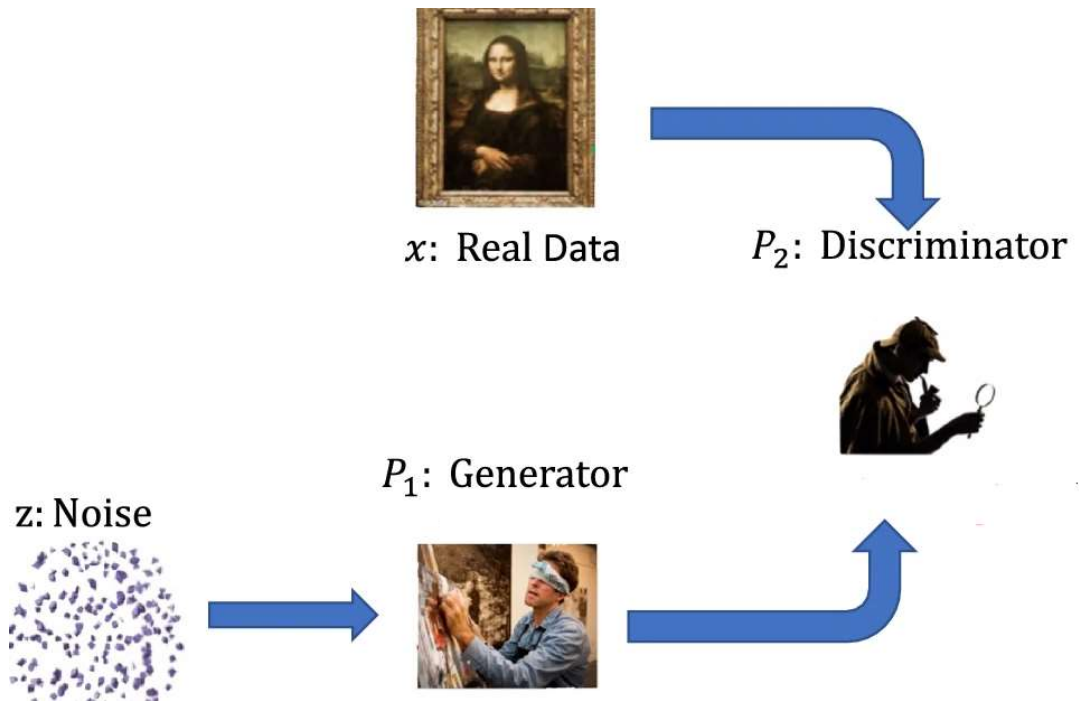


Adapted from [Sauer et al. (2022)]

## Class 2 of spatial probabilistic ML models: Generative like (GEN)

Translate the problem into learn the 'unknown' predictive distribution  $F_{X^*}^{X,Y}$  from the training data points

Based on the training data points  $(X_i, Y_i)_{i=1,\dots,n}$  learn,  $Y(X^*) \sim \text{Gauss}(\mu^*, \Sigma^*) \sim F_{X^*}^{X,Y}$



### Procedure:

1. Learn the **joint distribution**  $\mathcal{L}(Y, X)$  using  $P_1$
2. Predict at  $X^*$  by **conditioning**  
 $F_{X^*}^{X,Y} \sim \mathcal{L}(Y, X) | X = X^*$
3. **Generate** samples from  $F_{X^*}^{X,Y}$

Adversarial approach adapted from [Mohebbi Moghaddam et al. (2023)]

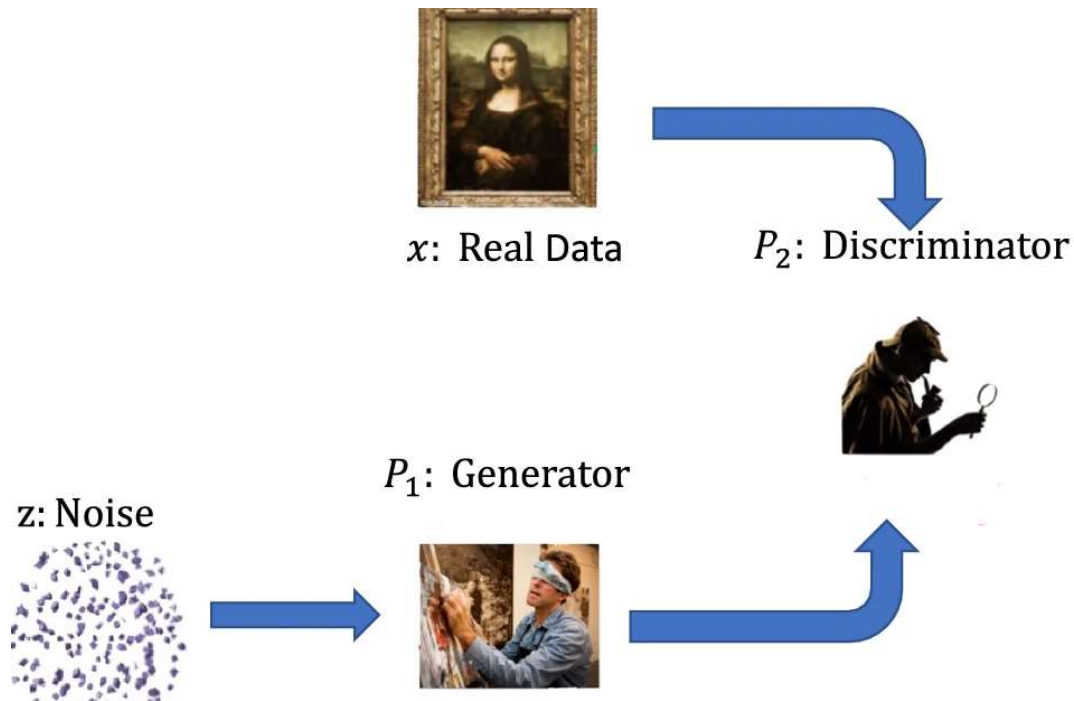
<https://arxiv.org/pdf/2106.06976>



## Class 2 of spatial probabilistic ML models: Generative like (GEN)

Translate the problem into learn the 'unknown' predictive distribution  $F_{X^*}^{X,Y}$  from the training data points

Based on the training data points  $(X_i, Y_i)_{i=1,\dots,n}$  learn,  $Y(X^*) \sim \text{Gauss}(\mu^*, \Sigma^*) \sim F_{X^*}^{X,Y}$



### Specificities of our problem:

□ Data are **tabular**

→ use of random forest RF instead of NN

[Watson et al., 2023]

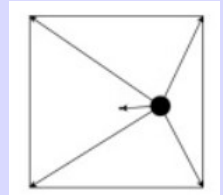
→ In this case,  $F^{X,Y}$  = mixture of 1d density distributions extracted from the RF leafs

□ **Spatial dependencies**

→ Introduce additional covariates corresponding to highly correlated spatial fields

→ Use of Euclidean Distance Fields

[Behrens et al., 2018]



> 16

## Class 3 of spatial probabilistic ML models: Conformal predictions (CF)

Translate the problem into assessing a **valid**  $PI^\alpha$   $\text{Prob}(Y^* \in PI^\alpha) \geq 1 - \alpha$  from the training data points

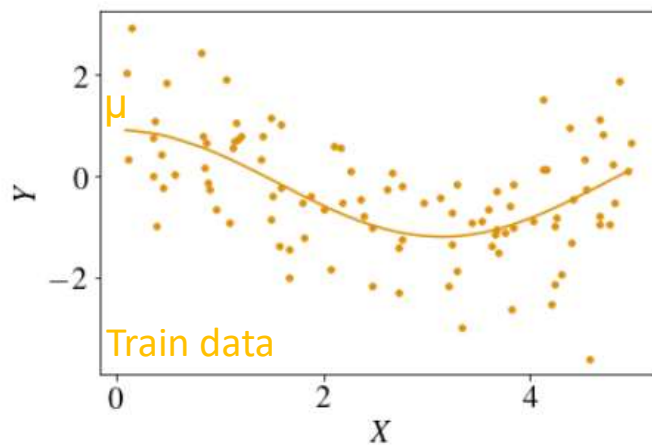
□ Use of **Split Conformal Prediction** (SCP) [Vovk et al. (2005); Papadopoulos et al. (2002), Lei et al. 2018]



## Class 3 of spatial probabilistic ML models: Conformal predictions (CF)

Translate the problem into assessing a **valid**  $PI^\alpha$   $\text{Prob}(Y^* \in PI^\alpha) \geq 1 - \alpha$  from the training data points

□ Use of **Split Conformal Prediction** (SCP) [Vovk et al. (2005); Papadopoulos et al. (2002), Lei et al. 2018]

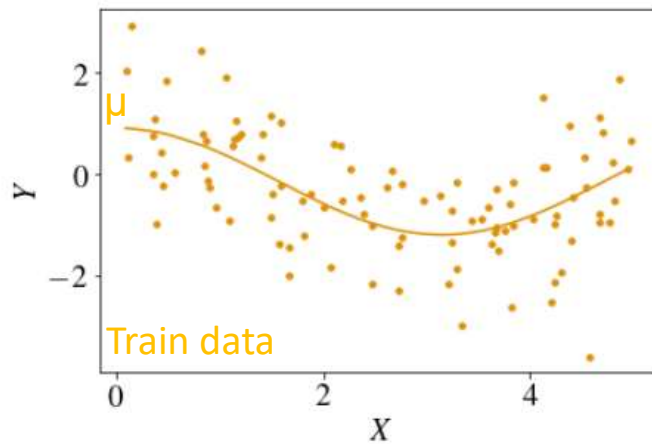


**Stage 1:** Estimate ML mean  $\mu$

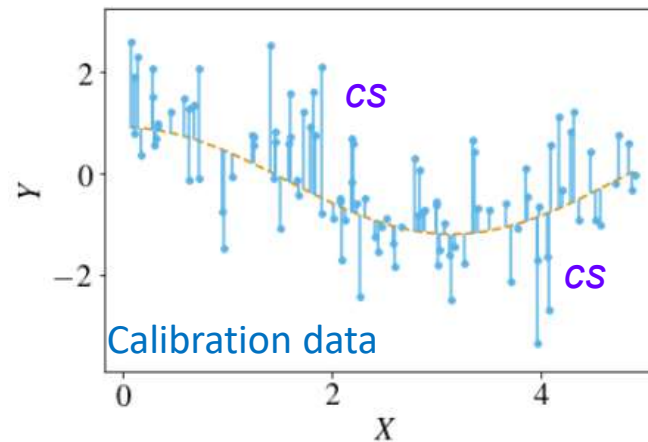
## Class 3 of spatial probabilistic ML models: Conformal predictions (CF)

Translate the problem into assessing a **valid**  $PI^\alpha$   $\text{Prob}(Y^* \in PI^\alpha) \geq 1 - \alpha$  from the training data points

□ Use of **Split Conformal Prediction** (SCP) [Vovk et al. (2005); Papadopoulos et al. (2002), Lei et al. 2018]



**Stage 1:** Estimate ML mean  $\mu$



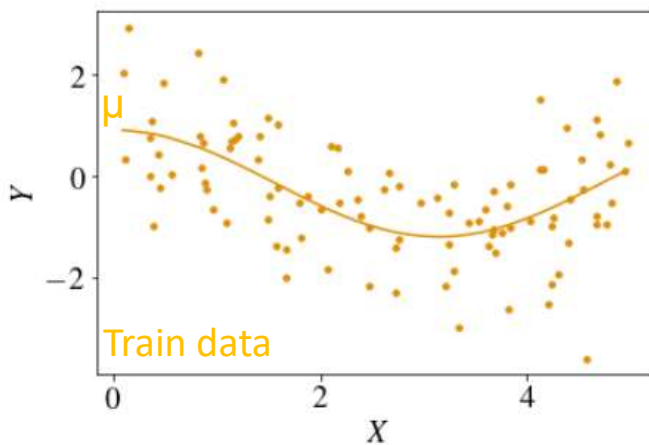
**Stage 2:** Estimate the non-conformity scores  $cs$  using  $\mu$



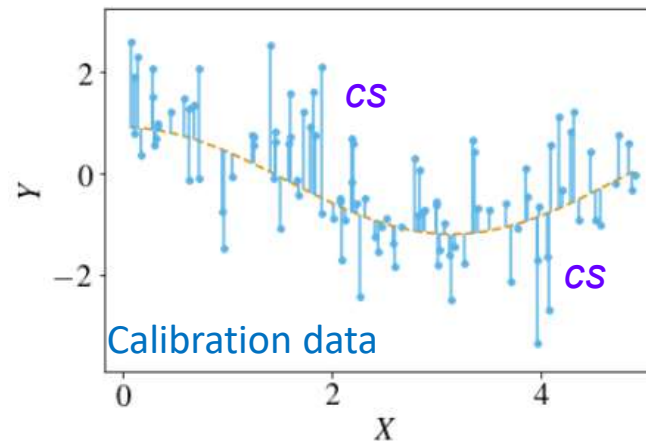
## Class 3 of spatial probabilistic ML models: Conformal predictions (CF)

Translate the problem into assessing a **valid**  $PI^\alpha$   $\text{Prob}(Y^* \in PI^\alpha) \geq 1 - \alpha$  from the training data points

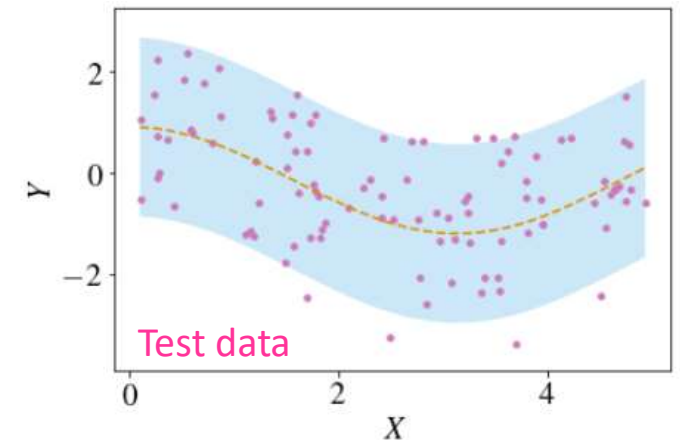
□ Use of **Split Conformal Prediction** (SCP) [Vovk et al. (2005); Papadopoulos et al. (2002), Lei et al. 2018]



**Stage 1:** Estimate ML mean  $\mu$



**Stage 2:** Estimate the non-conformity scores  $cs$  using  $\mu$



**Stage 3:** Compute the  $(1-\alpha)$  empirical quantile  $Q^{1-\alpha}(S)$  of  $S = \{cs\}_{cal} \cup \{+\infty\}$

$$\mathcal{I}((X_1, Y_1), \dots, (X_n, Y_n)) = \mathcal{I}((X_{\sigma(1)}, Y_{\sigma(1)}), \dots, (X_{\sigma(n)}, Y_{\sigma(n)}))$$

For any permutation  $\sigma$  of  $(1, \dots, n)$

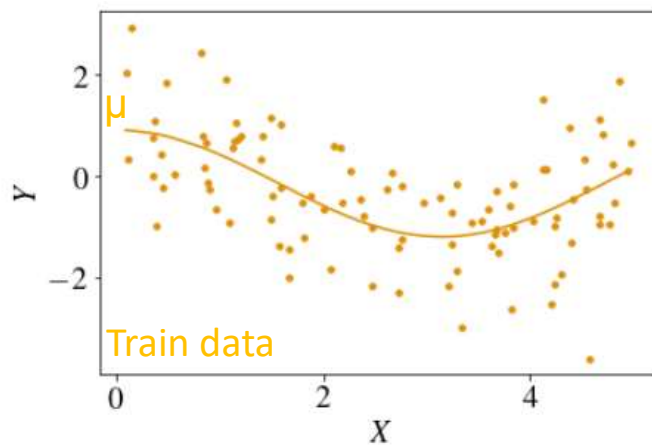


Calibration and test data need to be exchangeable!!

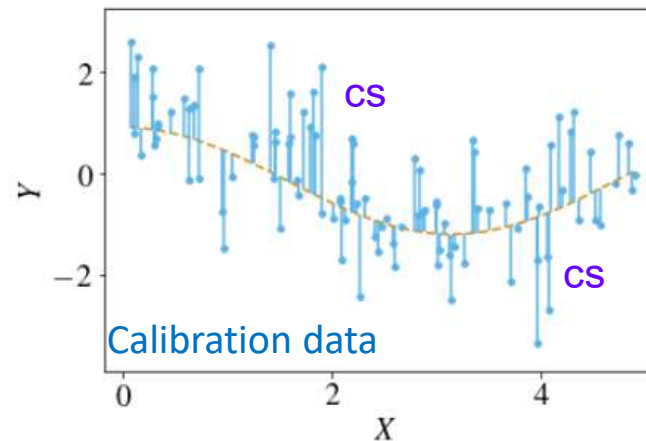
## Class 3 of spatial probabilistic ML models: Conformal predictions (CF)

Translate the problem into assessing a **valid**  $PI^\alpha$   $\text{Prob}(Y^* \in PI^\alpha) \geq 1 - \alpha$  from the training data points

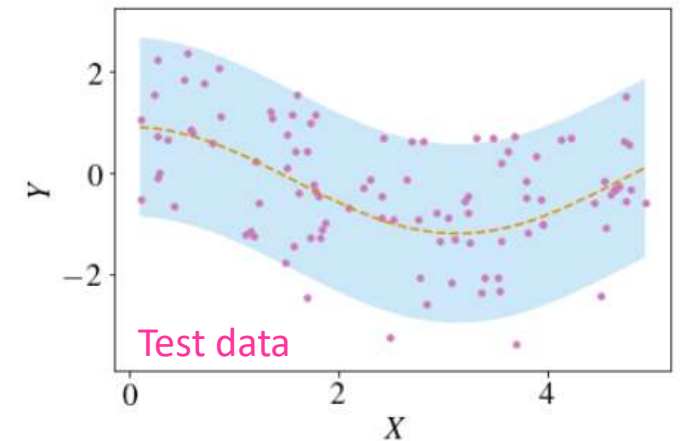
□ Use of **Split Conformal Prediction** (SCP) [Vovk et al. (2005); Papadopoulos et al. (2002), Lei et al. 2018]



**Stage 1:** Estimate ML mean  $\mu$



**Stage 2:** Estimate the non-conformity scores  $cs$



**Stage 3:** Compute the  $(1-\alpha)$  empirical quantile  $Q^{1-\alpha}(S)$  of  $S = \{cs\}_{cal} \cup \{+\infty\}$

□ Adaptation to the spatial context [Mao et al. (2020)]

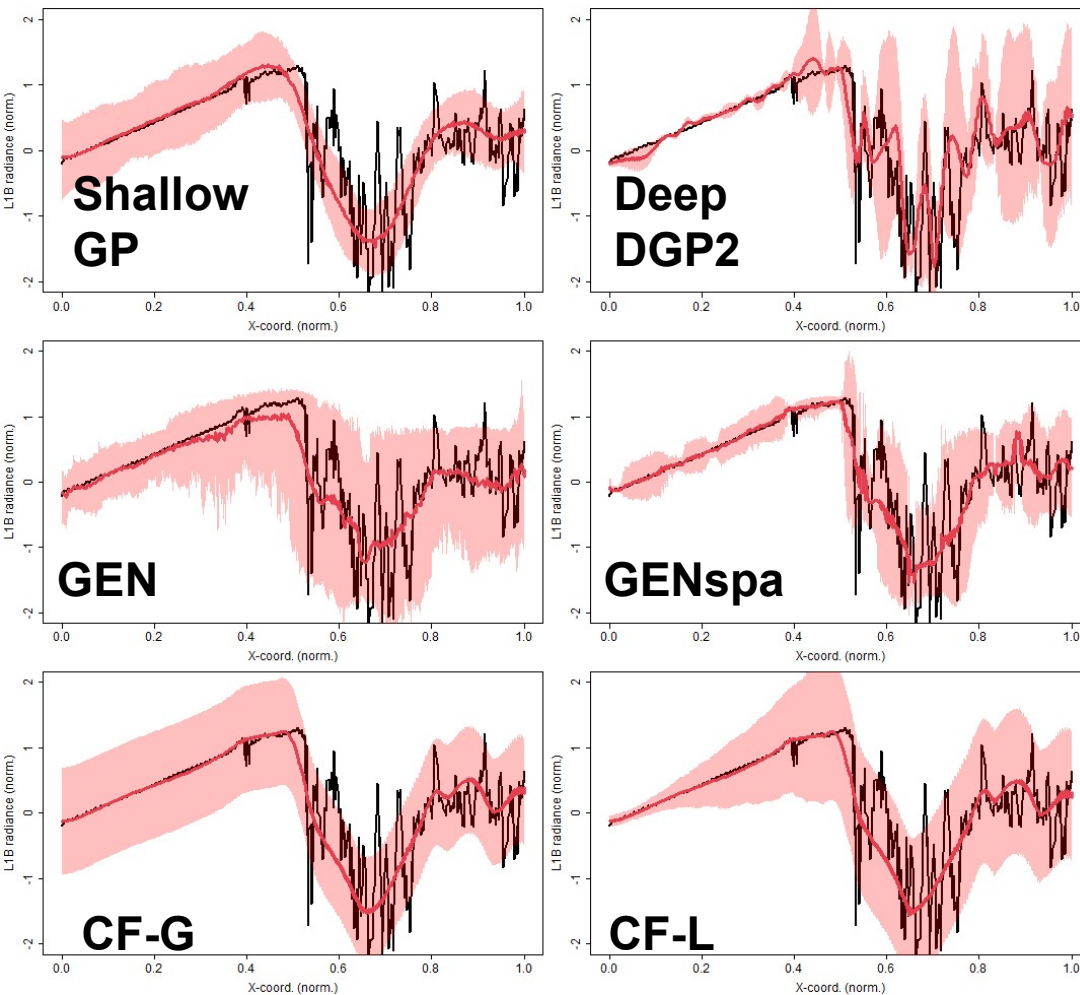
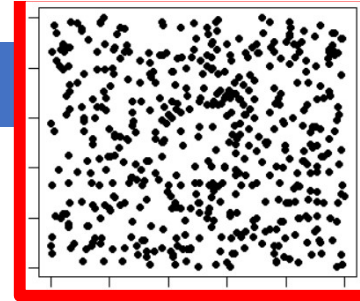
**Global**  $cs_i = \frac{|y_i - \mu(X_i)|}{\sigma(X_i)}$  where  $\mu, \sigma$  are given by a GP

**Local**

Same as **Global** but over a region around the prediction point determined via CV with maximisation of interval score



## An example of prediction – real case – **RANDOM, N=500**

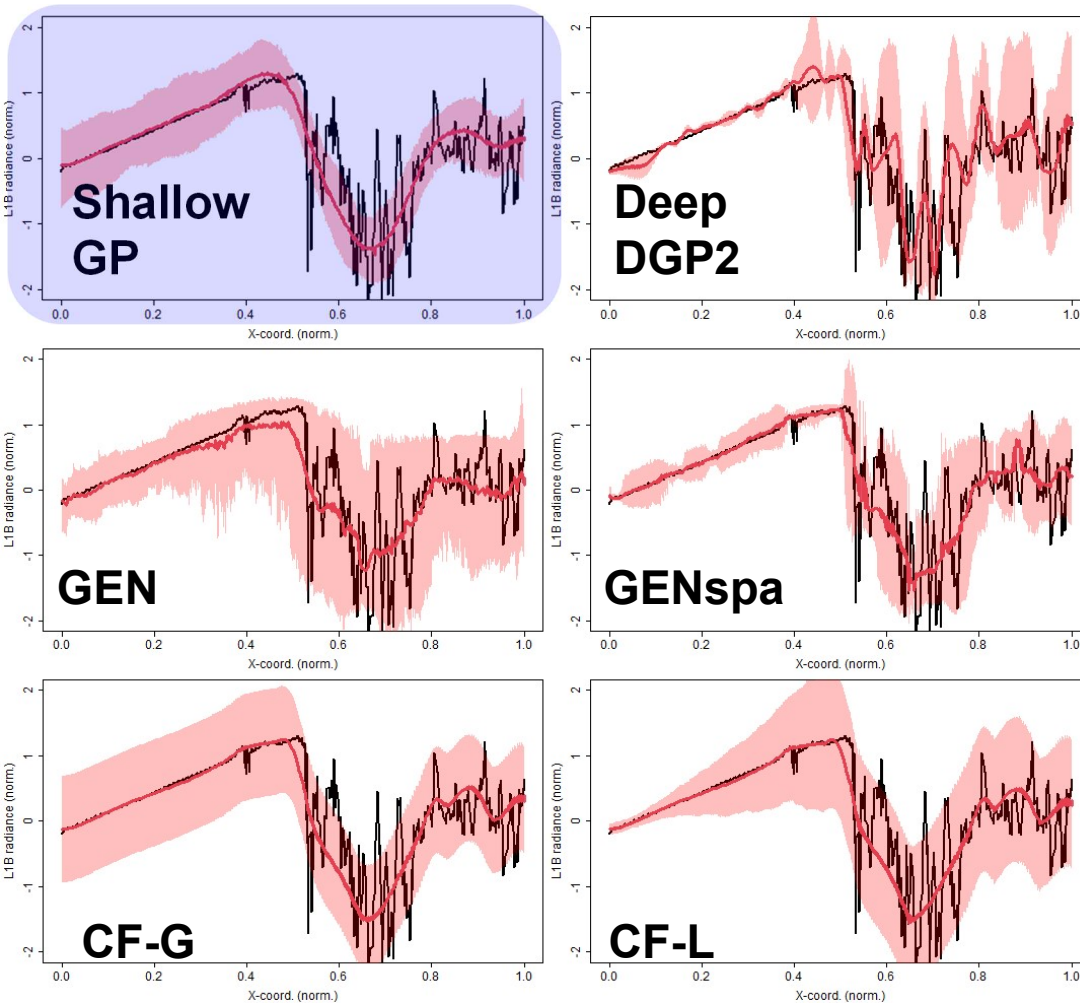
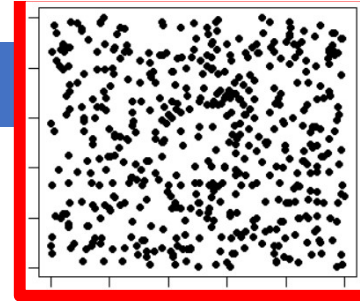


 90% unc. envelope     Mean

### Computation

- ☐ DGP, GEN: quantiles computed from a set of 500 stochastic simulations
- ☐ CF: direct use of the conformal predictions

## An example of prediction – real case – **RANDOM, N=500**



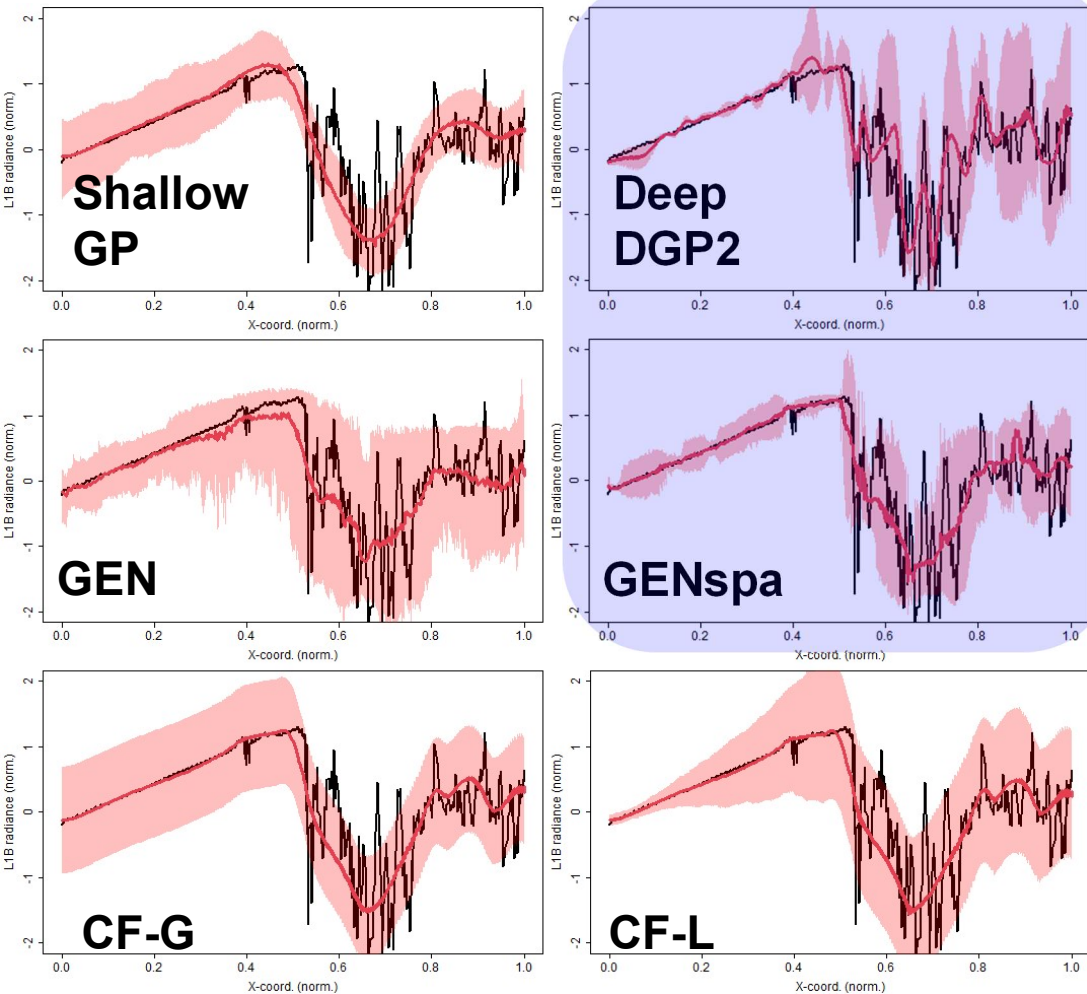
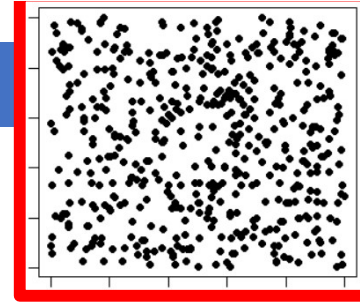
 90% unc. envelope     Mean

❑ Shallow GP captures **medium range** variations

### Computation

- ❑ DGP, GEN: quantiles computed from a set of 500 stochastic simulations
- ❑ CF: direct use of the conformal predictions

## An example of prediction – real case – **RANDOM, N=500**



 90% unc. envelope  Mean

❑ Shallow GP captures **medium range** variations

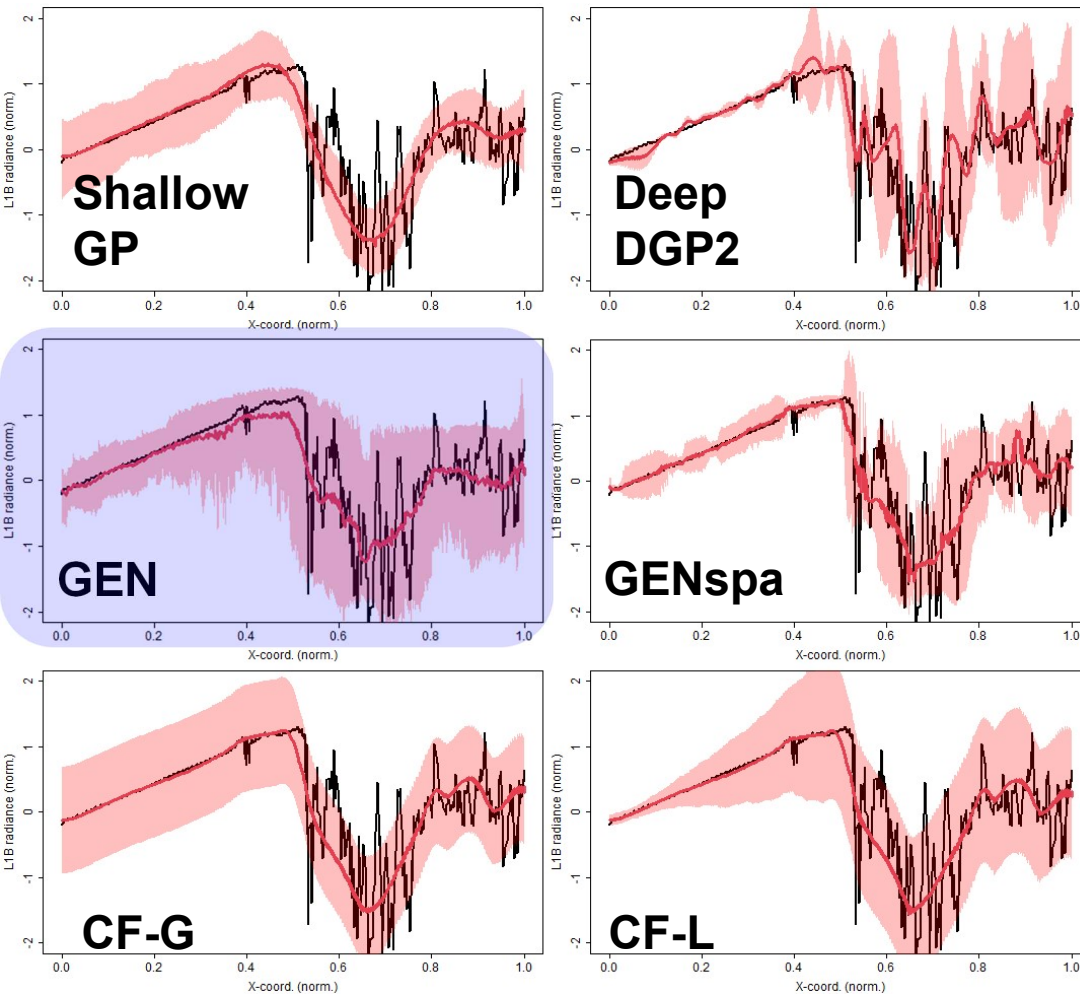
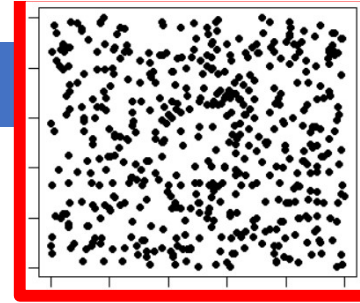
❑ DGP, GENspa capture variations of **multiple ranges of variation**

### Computation

- ❑ DGP, GEN: quantiles computed from a set of 500 stochastic simulations
- ❑ CF: direct use of the conformal predictions



## An example of prediction – real case – **RANDOM, N=500**



 90% unc. envelope     Mean

❑ Shallow GP captures **medium range** variations

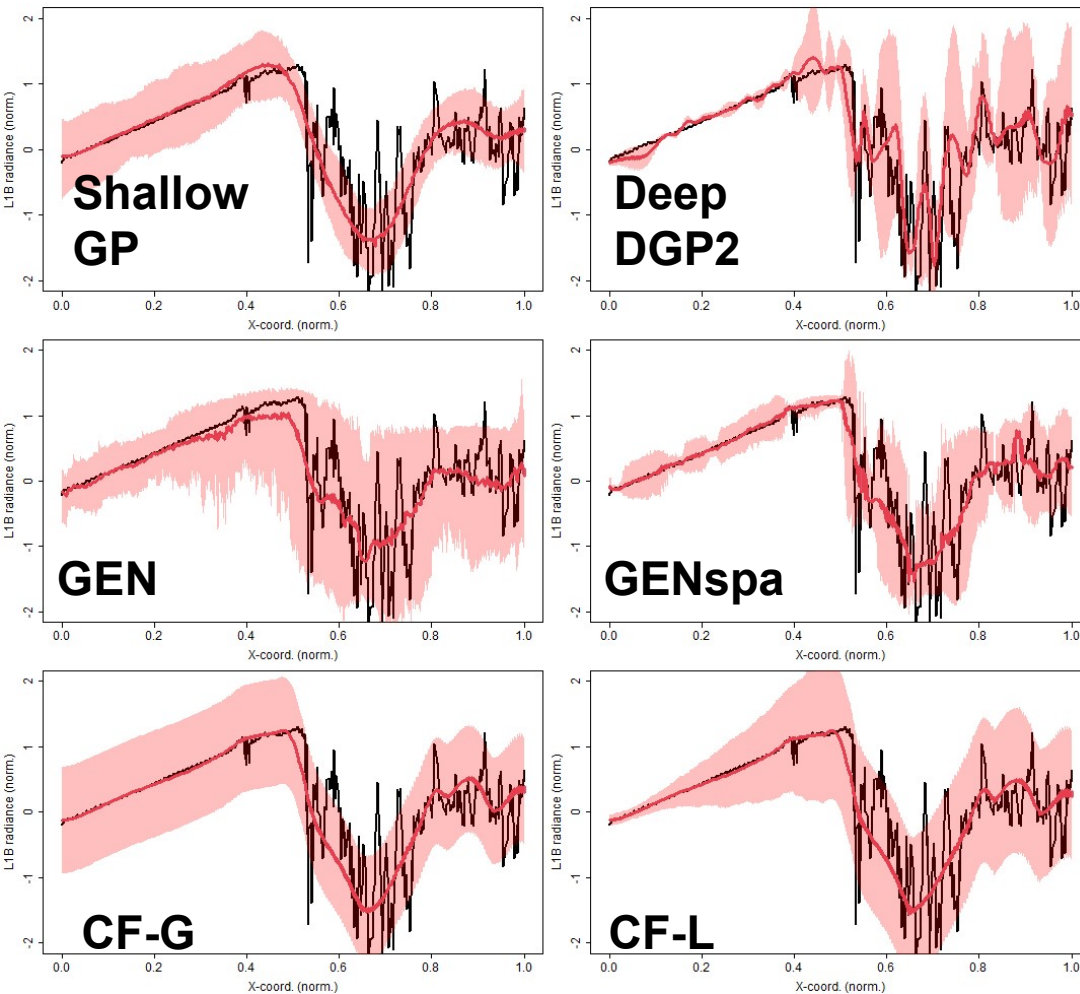
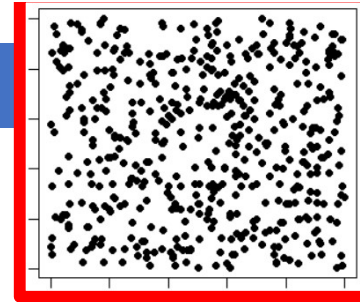
❑ DGP, GENspa capture variations of **multiple ranges of variation**

❑ GEN provides **too wide prediction intervals**

### Computation

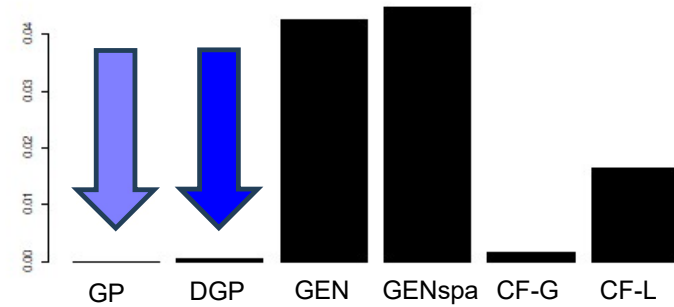
- ❑ DGP, GEN: quantiles computed from a set of 500 stochastic simulations
- ❑ CF: direct use of the conformal predictions

# An example of prediction – real case – **RANDOM, N=500**



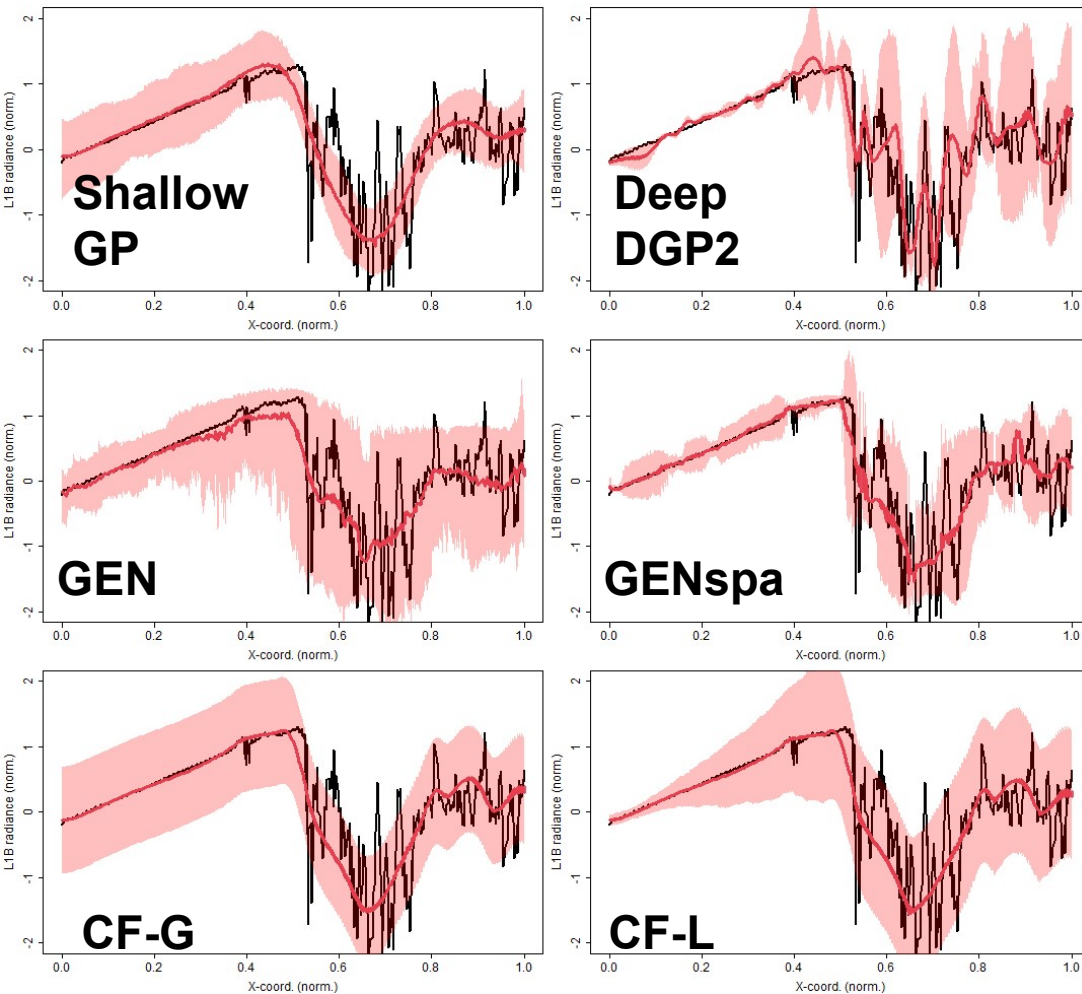
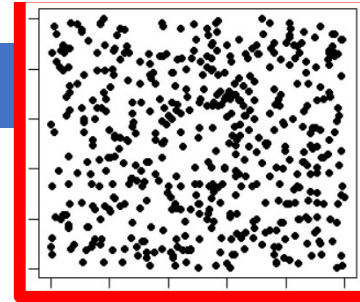
90% unc. envelope     Mean

Score – min(Score)



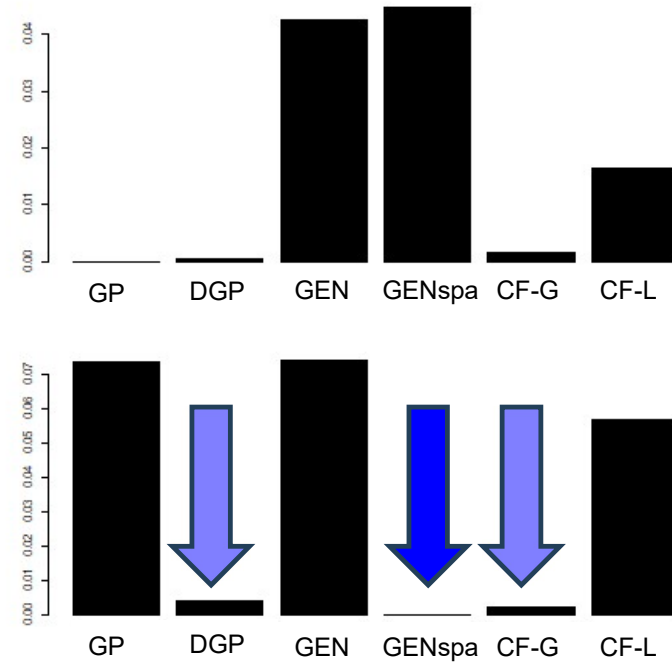
$Q^2$

# An example of prediction – real case – **RANDOM, N=500**



90% unc. envelope  Mean

Score – min(Score)

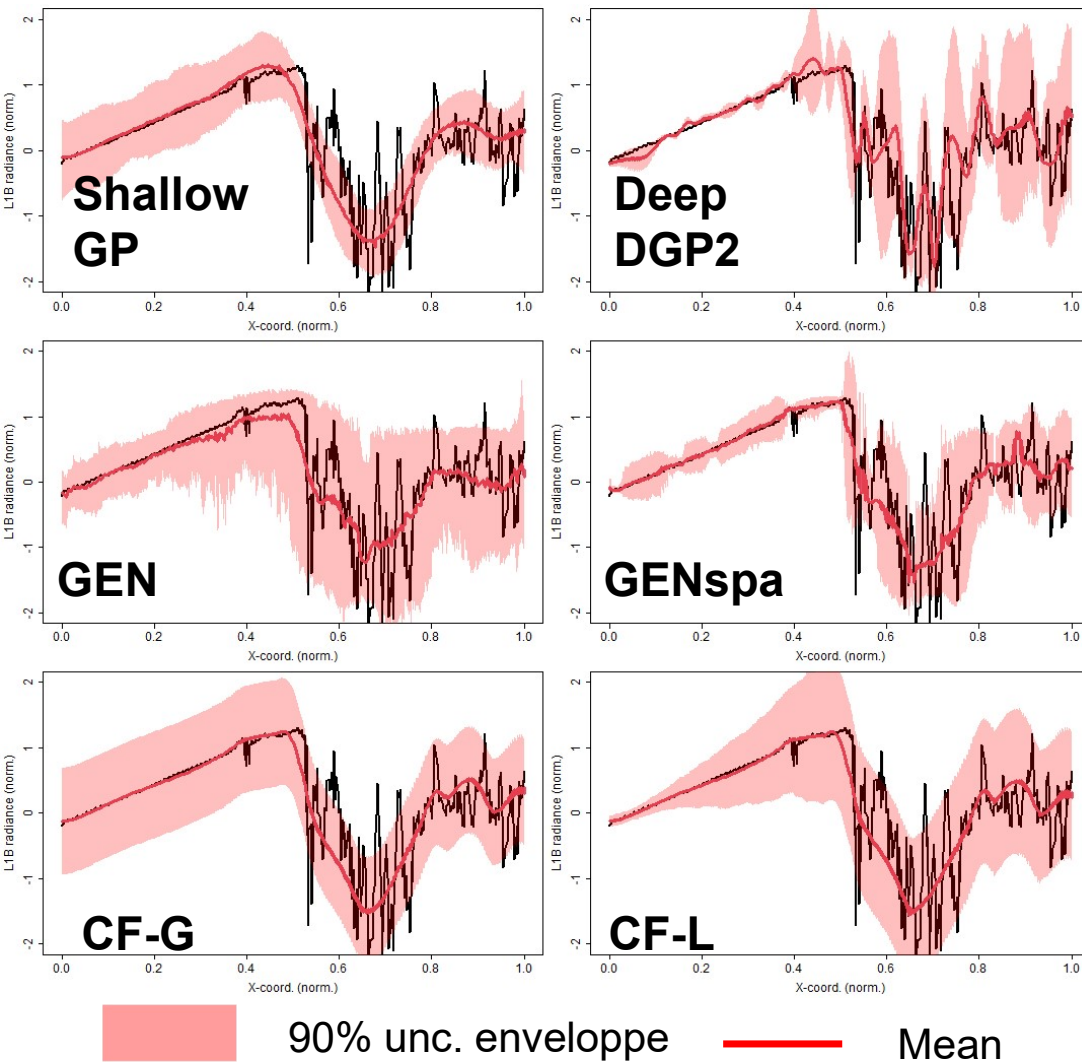
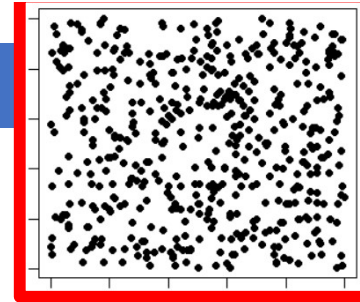


$Q^2$

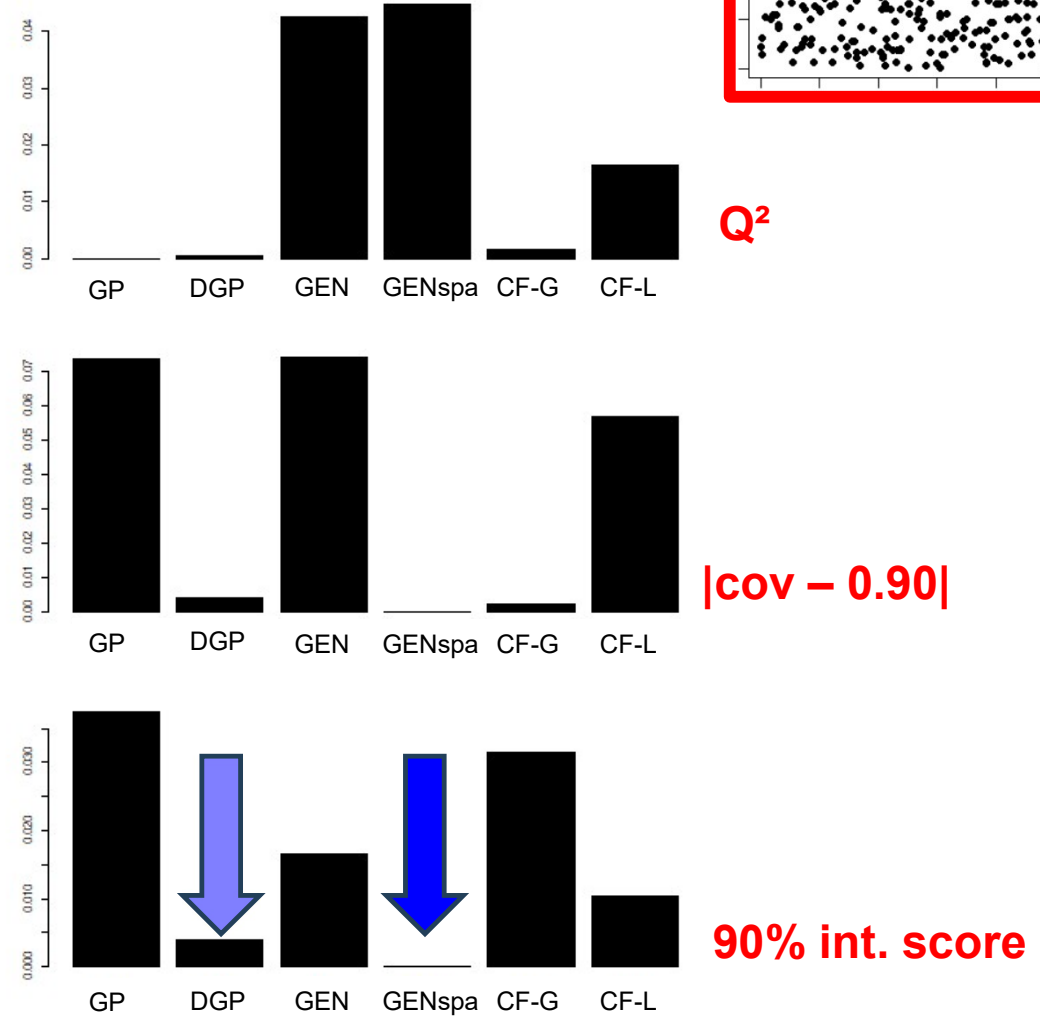
$|\text{cov} - 0.90|$



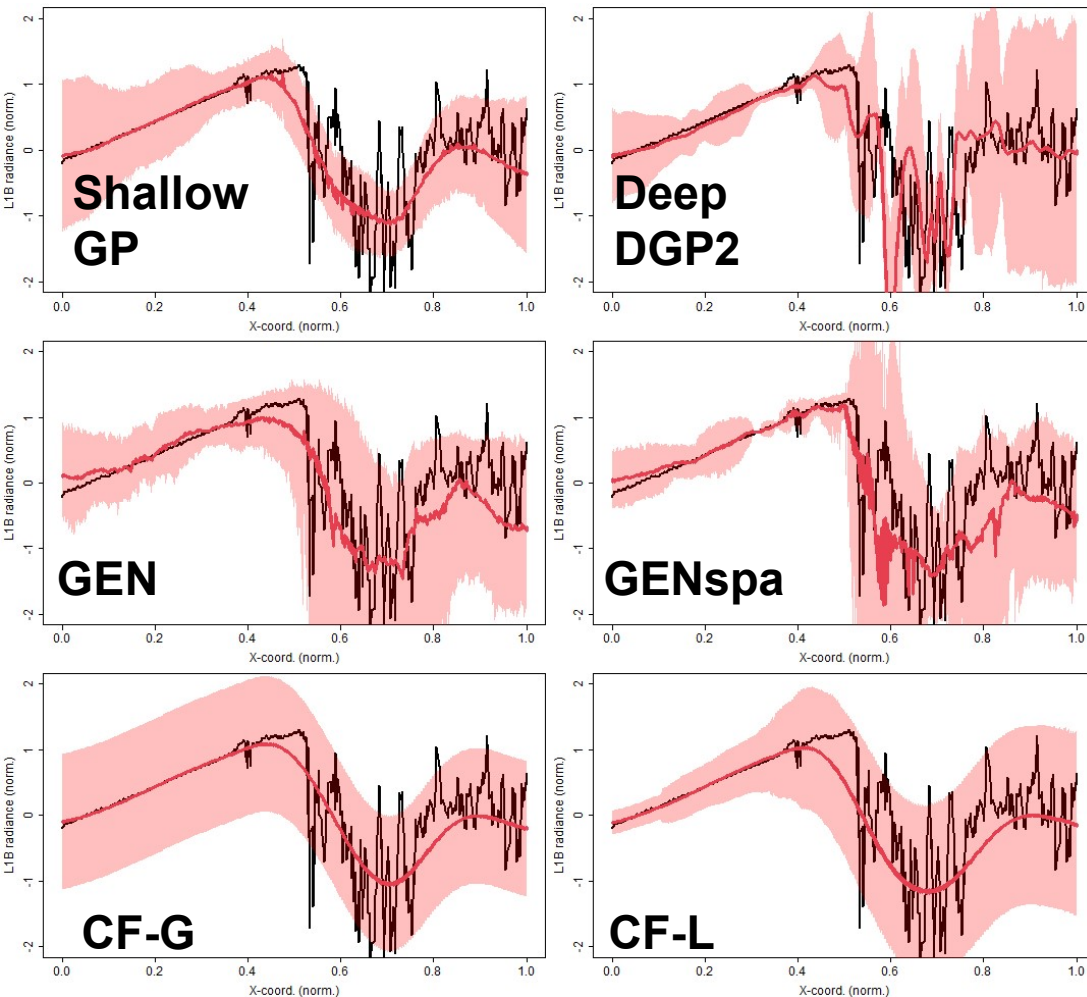
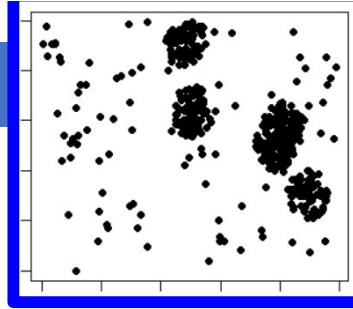
# An example of prediction – real case – **RANDOM, N=500**



Score – min(Score)

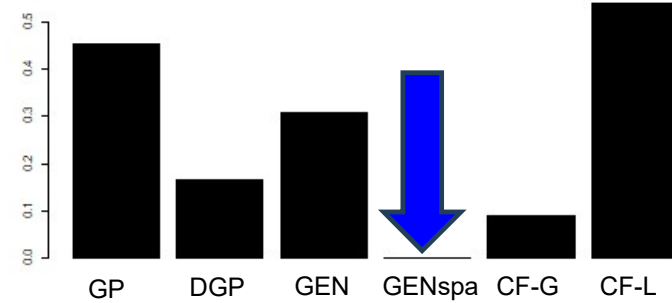


# An example of prediction – real case – CLUSTERED, N=500

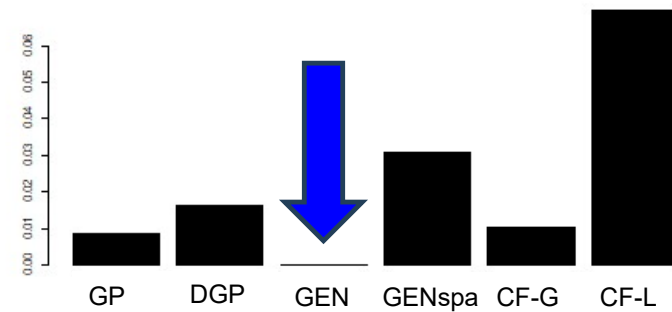


90% unc. envelope — Mean

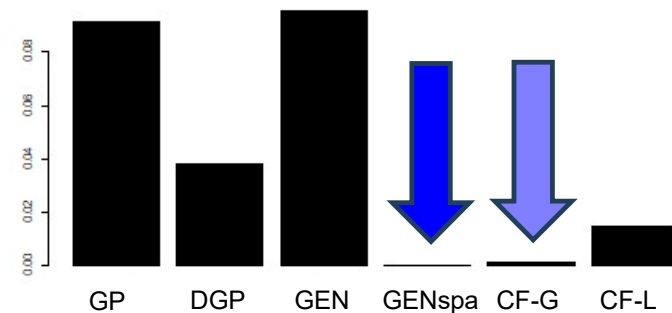
## Score – min(Score)



$Q^2$



|cov – 0.90|

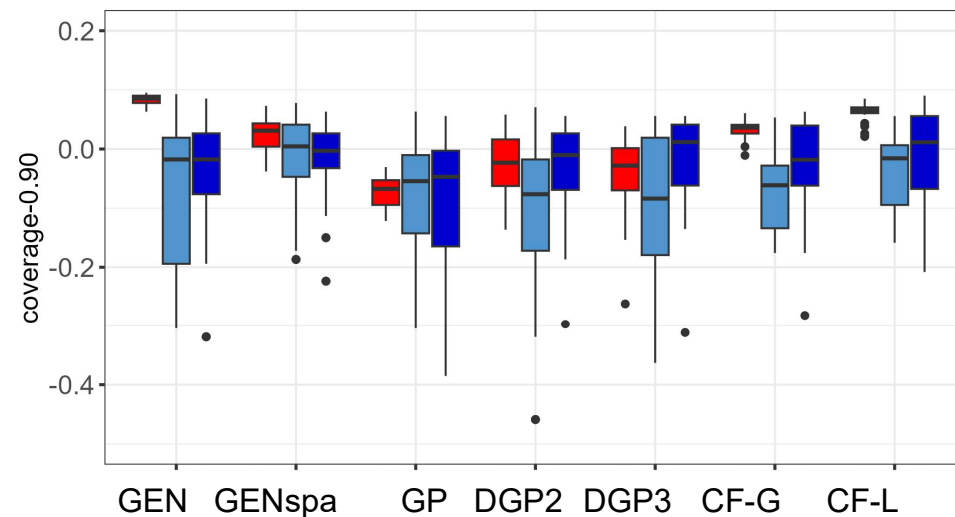
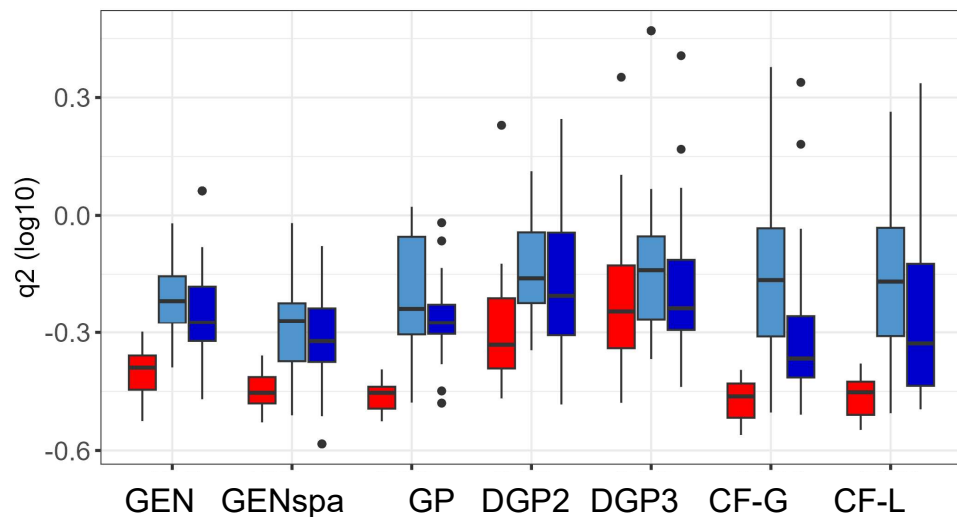


90% int. score

> 29

> 29

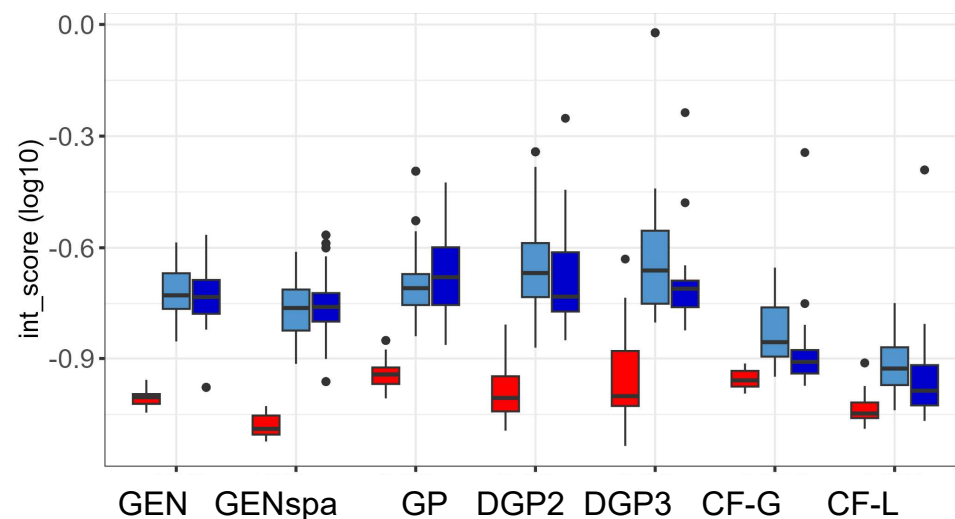
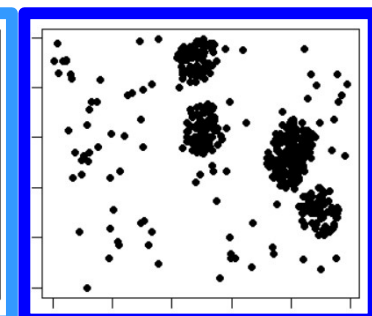
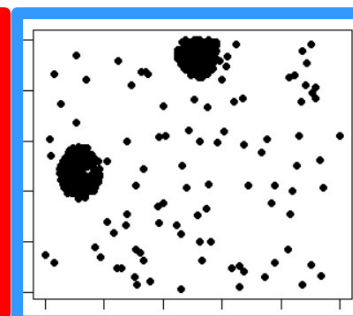
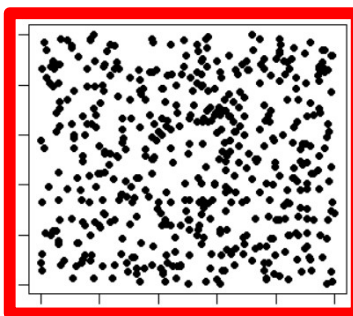
# Results of 25 repeated random experiments – real case



RANDOM, N=500

2 CLUSTERS

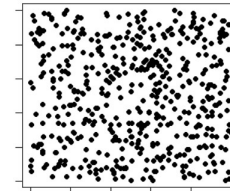
4 CLUSTERS





## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16



### Median value based on 25 repeated random experiments

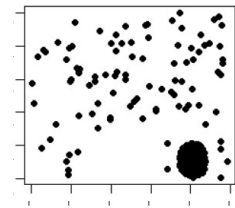
- ❑ Deep GP performs well for uncertainty-oriented scores
- ❑ Overall, GENspa is the best performing model

## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16
$1-Q^2$	0.67	<b>0.58</b>	0.77	0.82	0.81	0.59	0.60
0.9-Coverage	0.06	0.04	0.06	0.08	0.07	0.07	<b>0.03</b>
Interval score 90%	0.13	0.07	0.17	0.15	0.16	<b>0.05</b>	0.07
0.5-Coverage	0.20	0.19	0.21	0.22	0.24	<b>0.13</b>	<b>0.13</b>
Interval score 50%	0.25	<b>0.21</b>	0.28	0.27	0.29	0.24	0.23



1 cluster

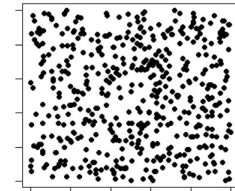


The clustering worsens performance:

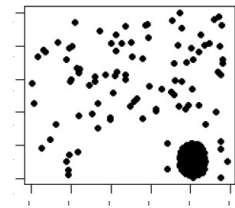
- ❑  $Q^2$  decreases by **~70%** (in average)
- ❑ Interval score for **moderate quantiles** increases by **120%** (in average)

## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16
$1-Q^2$	0.67	<b>0.58</b>	0.77	0.82	0.81	0.59	0.60
0.9-Coverage	0.06	0.04	0.06	0.08	0.07	0.07	<b>0.03</b>
Interval score 90%	0.13	0.07	0.17	0.15	0.16	<b>0.05</b>	0.07
0.5-Coverage	0.20	0.19	0.21	0.22	0.24	<b>0.13</b>	<b>0.13</b>
Interval score 50%	0.25	<b>0.21</b>	0.28	0.27	0.29	0.24	0.23



**1 cluster**

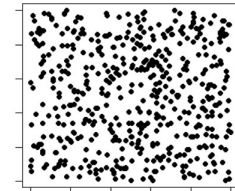


❑ Shallow or Deep GP performance worsens due to clustering

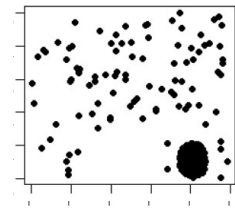


## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16
$1-Q^2$	0.67	<b>0.58</b>	0.77	0.82	0.81	0.59	0.60
0.9-Coverage	0.06	0.04	0.06	0.08	0.07	0.07	<b>0.03</b>
Interval score 90%	0.13	0.07	0.17	0.15	0.16	<b>0.05</b>	0.07
0.5-Coverage	0.20	0.19	0.21	0.22	0.24	<b>0.13</b>	<b>0.13</b>
Interval score 50%	0.25	<b>0.21</b>	0.28	0.27	0.29	0.24	0.23



1 cluster



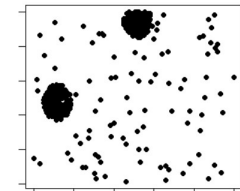
- ❑ CF performs relatively well
- ❑ Overall, GENspa is the best performing model

## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16
$1-Q^2$	0.60	<b>0.54</b>	0.58	0.69	0.73	0.68	0.68
0.9-Coverage	0.07	0.05	0.06	0.08	0.08	0.06	<b>0.03</b>
Interval score 90%	0.17	<b>0.07</b>	0.17	0.12	0.13	0.09	0.09
0.5-Coverage	0.19	0.17	0.20	0.21	0.22	0.14	<b>0.12</b>
Interval score 50%	0.22	<b>0.19</b>	0.24	0.28	0.30	0.27	0.24



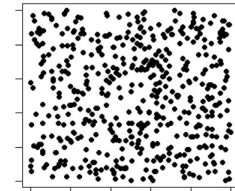
**2 clusters**



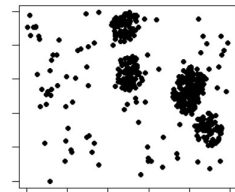
- ❑ Same result for GENspa and CF
- ❑ 2 clusters → more distributed information → GP slightly performs better

## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16
$1-Q^2$	0.53	0.48	0.53	0.62	0.58	<b>0.43</b>	0.47
0.9-Coverage	0.05	<b>0.03</b>	0.06	0.04	0.04	0.05	0.06
Interval score 90%	0.10	<b>0.05</b>	0.10	0.09	0.10	0.08	0.11
0.5-Coverage	0.19	0.17	0.21	0.19	0.19	0.12	<b>0.10</b>
Interval score 50%	0.21	<b>0.18</b>	0.25	0.22	0.22	0.22	0.20



4 clusters

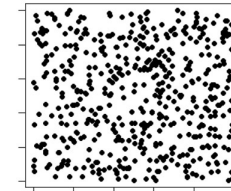


- ❑ Same conclusion as with 2 clusters
- ❑ 4 clusters → Even more distributed info. → some improvement of DGP

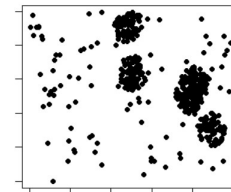


# Synthesis – real case – median over 25 random experiments

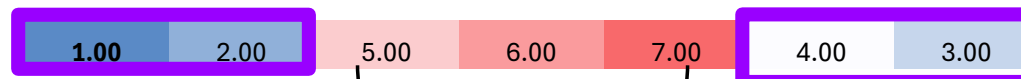
	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16



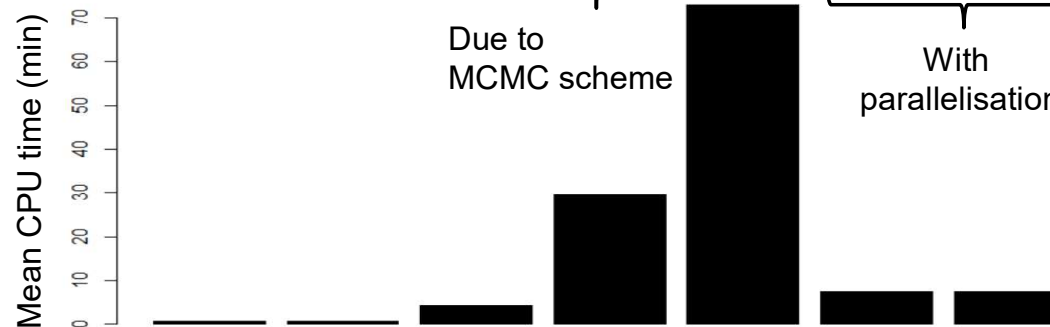
$1-Q^2$	0.53	0.48	0.53	0.62	0.58	<b>0.43</b>	0.47
0.9-Coverage	0.05	<b>0.03</b>	0.06	0.04	0.04	0.05	0.06
Interval score 90%	0.10	<b>0.05</b>	0.10	0.09	0.10	0.08	0.11
0.5-Coverage	0.19	0.17	0.21	0.19	0.19	0.12	<b>0.10</b>
Interval score 50%	0.21	<b>0.18</b>	0.25	0.22	0.22	0.22	0.20



CPU time\*

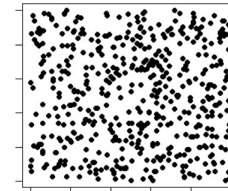


\*ranking based on 25 tests

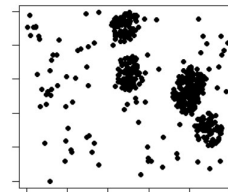


## Synthesis – real case – median over 25 random experiments

	GEN	GENspa	GP	DGP2	DGP3	CF-G	CF-L
$1-Q^2$	0.41	0.35	0.35	0.47	0.57	<b>0.34</b>	0.35
0.9-Coverage	0.09	<b>0.03</b>	0.07	0.04	0.04	0.04	0.07
Interval score 90%	0.10	<b>0.08</b>	0.11	0.10	0.10	0.11	0.09
0.5-Coverage	0.17	0.08	<b>0.02</b>	0.07	0.05	0.10	0.15
Interval score 50%	<b>0.08</b>	0.09	0.13	0.11	0.12	0.20	0.16



$1-Q^2$	0.53	0.48	0.53	0.62	0.58	<b>0.43</b>	0.47
0.9-Coverage	0.05	<b>0.03</b>	0.06	0.04	0.04	0.05	0.06
Interval score 90%	0.10	<b>0.05</b>	0.10	0.09	0.10	0.08	0.11
0.5-Coverage	0.19	0.17	0.21	0.19	0.19	0.12	<b>0.10</b>
Interval score 50%	0.21	<b>0.18</b>	0.25	0.22	0.22	0.22	0.20



CPU time*	<b>1.00</b>	2.00	5.00	6.00	7.00	4.00	3.00
Impl. Effort**	<b>1.00</b>	1.00	5.00	7.00	7.00	2.00	4.00

Rapid convergence,  
few hyperparameters

Careful  
convergence analysis

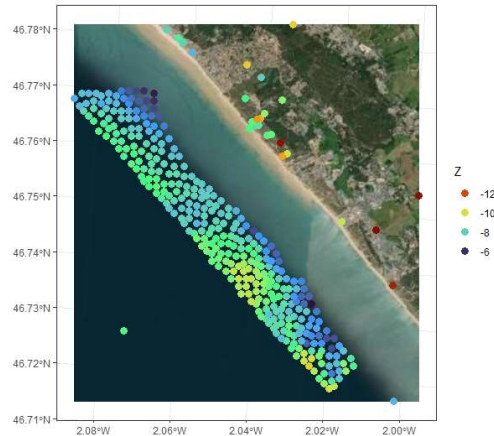
Size of neighbour region  
difficult to assess

\*\*ranking based  
on my feedback

## Summary

- ❑ **Complex sample distributions** (cluster, sparse) result in **performance decline** (prediction accuracy AND uncertainty)
- ❑ **Deep Gaussian Process** performs well for **random settings** (coverage, interval score) but at the CPU time cost, + convergence checking
- ❑ **Conformal predictions** have an **intermediate performance**; no/slight improvement of the local version
- ❑ **Generative model** is **robust to the presence of clusters**, but need adequate modelling of spatial dependence
- ❑ **Results checked** also by varying the size of the clusters, number of samples, number of samples outside the clustered region, the type of benchmark cases...
- ❑ **Next step?** How to do when the ground truth is not available  
→ **cross validation for spatial data?**

# Open question: validity of a standard 10-fold random cross validation?



## ARTICLE

<https://doi.org/10.1016/j.ecolmodel.2021.109690>

OPEN

Spatial validation reveals poor predictive performance of large-scale ecological mapping models

Pierre Ploton<sup>1</sup>, Frédéric Mortier<sup>2,3</sup>, Maxime Réjou-Mé<sup>4</sup>, Vivien Rossi<sup>5</sup>, Carsten Dormann<sup>6</sup>, Guillaume Cornu<sup>7</sup>, Alexei Lyapunov<sup>8</sup>, Sylvie Gourlet-Fleury<sup>2,3</sup> & Raphaël Pélissier<sup>1</sup>

## COMMENT

<https://doi.org/10.1016/j.ecolmodel.2021.109690>

OPEN

Machine learning-based global maps of ecological variables and the challenge of assessing them

Hanna Meyer<sup>1</sup> & Edzer Pebesma<sup>2</sup>

Ecological Modelling 457 (2021) 109690

Contents lists available at ScienceDirect

Ecological Modelling

journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



Short communication

Spatial cross-validation is not the right way to evaluate map accuracy

Alexandre M.J.-C. Wadoux<sup>a,\*</sup>, Gerard B.M. Heuvelink<sup>b</sup>, Sytze de Bruin<sup>c</sup>, Dick J. Brus<sup>d</sup>



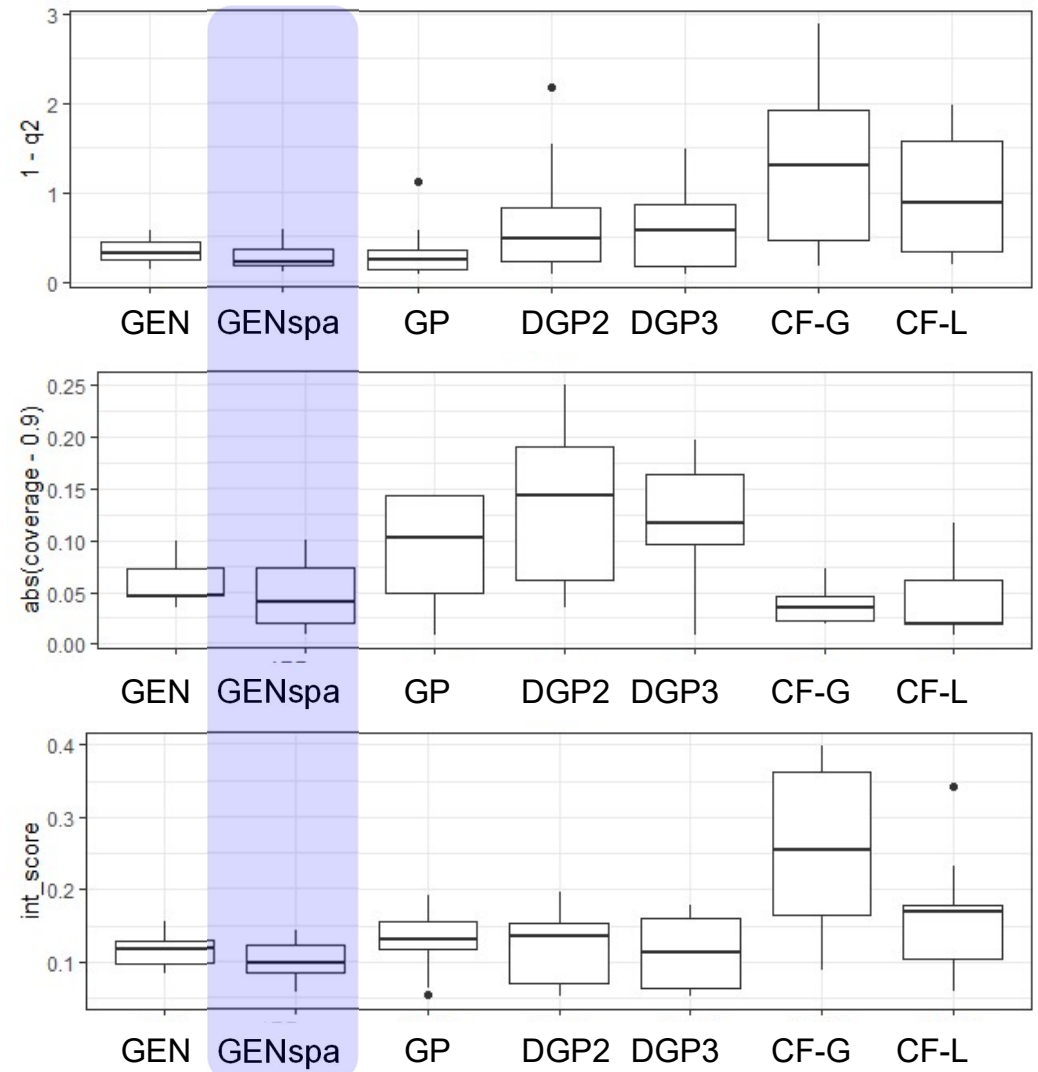
Contents lists available at ScienceDirect

Spatial Statistics

journal homepage: [www.elsevier.com/locate/spasta](http://www.elsevier.com/locate/spasta)

Automatic cross-validation in structured models: Is it time to leave out leave-one-out?

Ariz Adin<sup>a,1</sup>, Elias Teixeira Krainski<sup>2</sup>, Amanda Lenzi<sup>3</sup>, Zhedong Liu<sup>4</sup>, Joaquín Martínez-Minaya<sup>5</sup>, Håvard Rue<sup>2</sup>





# Thank you for your attention!

# Merci pour votre attention!

We acknowledge BRGM for  
providing **SAPHIR computing  
and storage resources**

We acknowledge financial funding by  
ANR-HOUSES  
(grant number: **ANR-22-CE56-0006**)  
<https://anrhouses.github.io/>



**anr**®  
agence nationale  
de la recherche

# References

- **Behrens, T., Schmidt, K., Viscarra Rossel, R. A., Gries, P., Scholten, T., & MacMillan, R. A.** (2018). Spatial modelling with Euclidean distance fields and machine learning. *European journal of soil science*, 69(5), 757-770.
- **Gneiting, T., & Raftery, A. E.** (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477), 359-378.
- **Sauer, A., Gramacy, R. B., & Higdon, D.** (2023). “Active Learning for Deep Gaussian Process Surrogates.” *Technometrics* 65 (1): 4–18.
- **Watson, D. S., Blesch, K., Kapar, J., & Wright, M. N.** (2023). Adversarial random forests for density estimation and generative modeling. In *International Conference on Artificial Intelligence and Statistics* (pp. 5357-5375). PMLR.
- **Wikle, C. K., & Zammit-Mangion, A.** (2023). Statistical deep learning for spatial and spatiotemporal data. *Annual Review of Statistics and Its Application*, 10(1), 247-270.
- **Williams, C. K., & Rasmussen, C. E.** (2006). *Gaussian processes for machine learning* (Vol. 2, No. 3, p. 4). Cambridge, MA: MIT press.
- **Vovk, V., Gammerman, A., & Shafer, G.** (2005). *Algorithmic learning in a random world*. Boston, MA: Springer US.
- **Zammit-Mangion, A., Ng, T. L. J., Vu, Q., & Filippone, M.** (2022). Deep compositional spatial models. *Journal of the American Statistical Association*, 117(540), 1787-1808.

# Appendices

1. Fit unsupervised random forest (Shi and Horvath, 2006): First, permute feature values in the given dataset  $\mathbf{X}$  randomly across instances to create naive synthetic dataset  $\tilde{\mathbf{X}}$ . Then, fit a random forest  $\hat{f}^0$  to distinguish instances from  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  (labeled accordingly), where splits in the forest's trees pick up the data's dependency structure.
2. If the accuracy of  $\hat{f}^0$  is above 50%, new synthetic data is sampled from the leaves of forest  $\hat{f}^0$  (generator step) and a new random forest  $\hat{f}^1$  is fit to classify real and synthetic data (discriminator step).
3. Data generation and discrimination is continued for  $k$  iterations until the accuracy of  $\hat{f}^k$  drops down to 50% or below. This indicates that the algorithm has converged, implying that all feature dependencies have been learned and features are mutually independent in the leaves.
4. FORDE step (density estimation): The estimated joint density  $\hat{p}_{\text{ARF}}$  can – thanks to the mutual independence assumption of features within the leaves – be formulated as a mixture of products  $\hat{p}_l$  of univariate densities  $\hat{p}_{lj}$  for leaf  $l$  and feature  $j$ , which can be estimated with any arbitrary univariate density estimator within the random forest's leaves, weighted by the share of real data  $\pi_l$  that falls into  $l$ :

$$\hat{p}_{\text{ARF}}(\mathbf{x}) = \sum_l \pi_l \hat{p}_l(\mathbf{x}) = \sum_l \pi_l \prod_j \hat{p}_{lj}(x_j).$$

5. FORGE step (data generation): Synthetic data is generated by drawing a leaf  $l$  from the random forest with probability  $\pi_l$  and then sampling from the estimated univariate densities  $\hat{p}_{lj}$  within that leaf.
- Once  $\hat{p}_{\text{ARF}}$  is estimated, ARF allows us to derive estimated conditional densities  $\hat{p}_{\text{ARF}}(x_j | \mathbf{X}_C = \mathbf{x}_C)$  for fixed values  $\mathbf{x}_C$  with arbitrary conditioning sets  $C$  without the need of refitting the ARF:

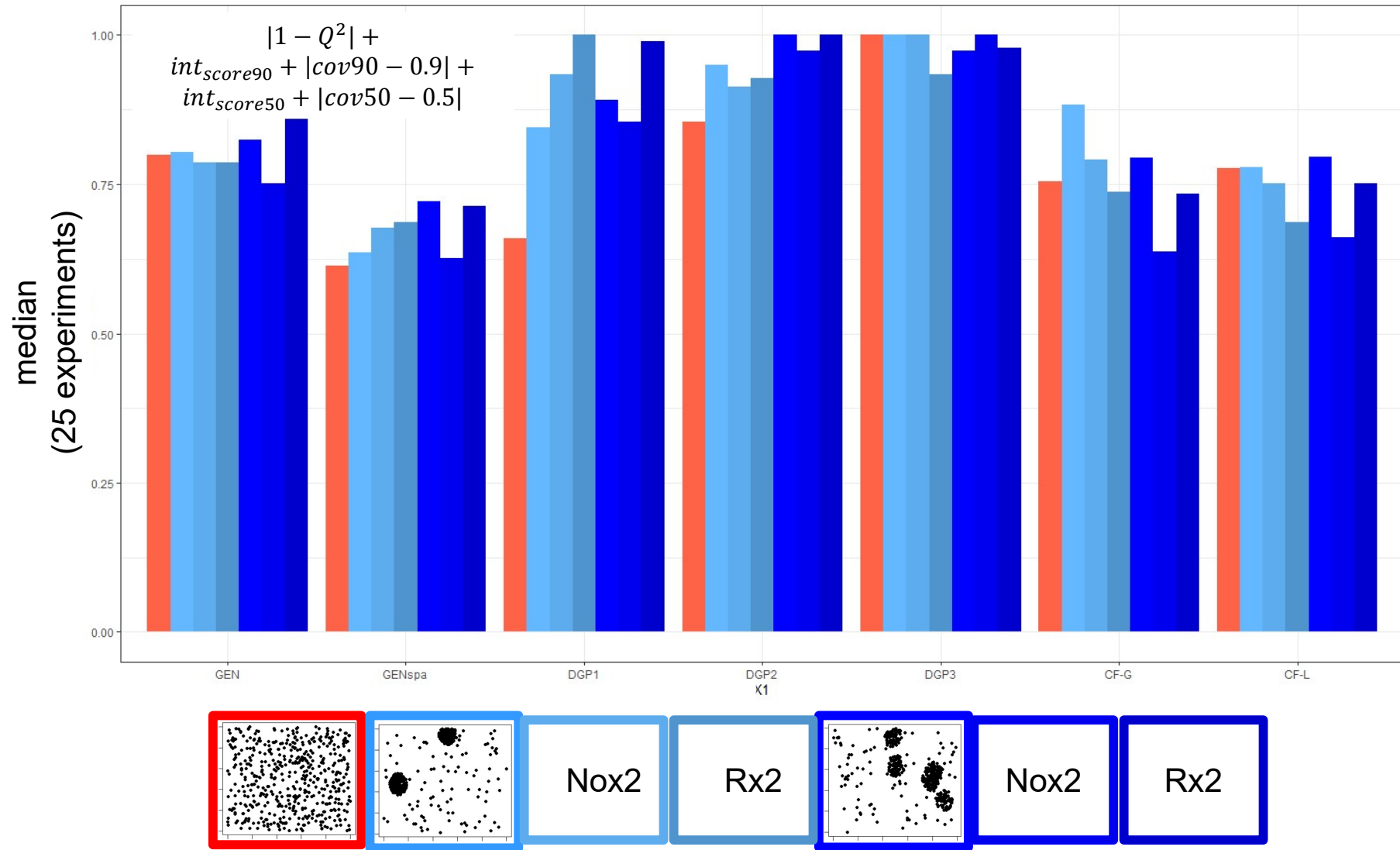
$$\hat{p}_{\text{ARF}}(x_j | \mathbf{X}_C = \mathbf{x}_C) = \sum_l \pi'_l \hat{p}_{lj}(x_j)$$

with updated weights  $\pi'_l := \pi_l \frac{\hat{p}_l(\mathbf{x}_C)}{\hat{p}_{\text{ARF}}(\mathbf{x}_C)}$ .

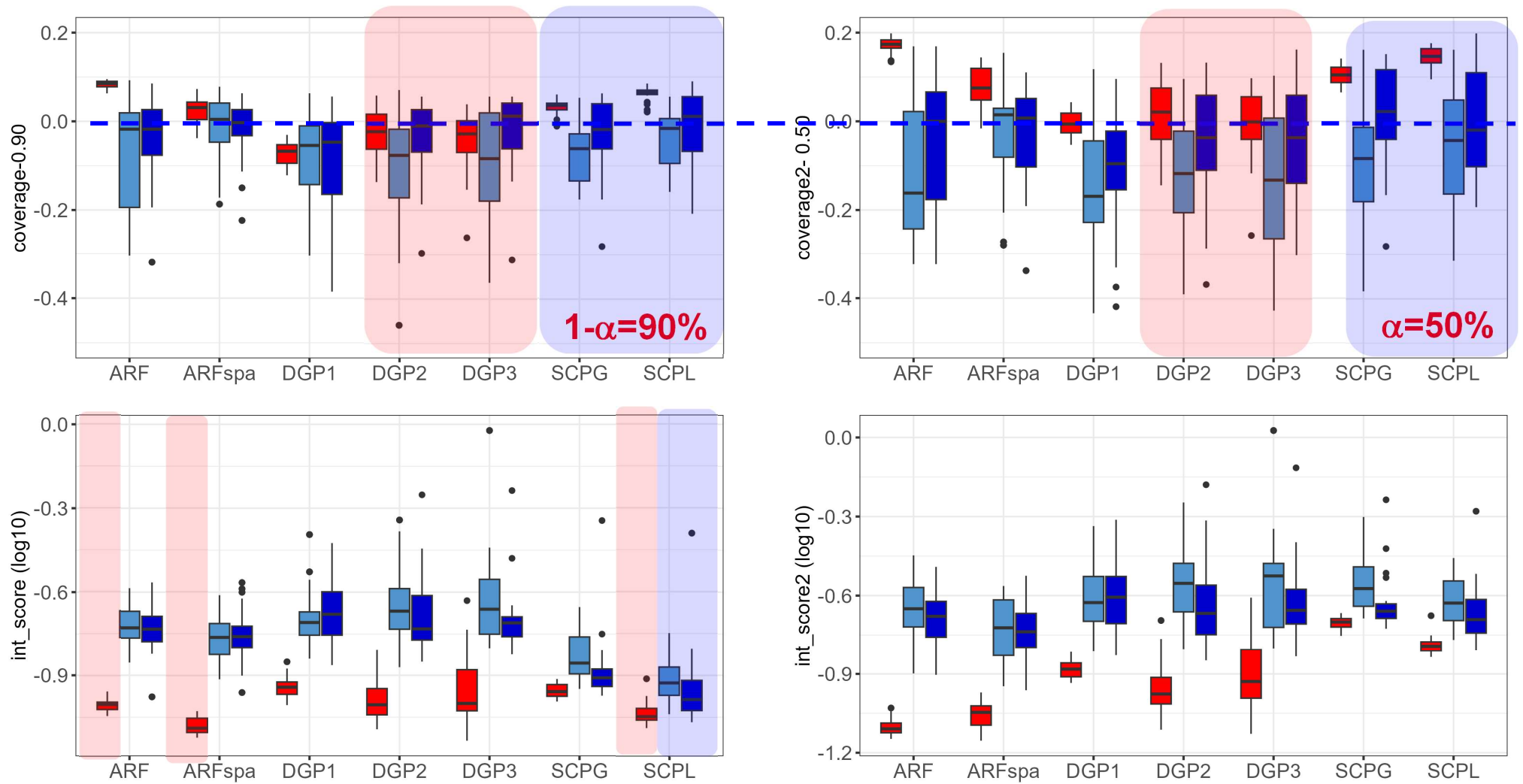
Watson et al. (2023); Blesch et al. (2025)



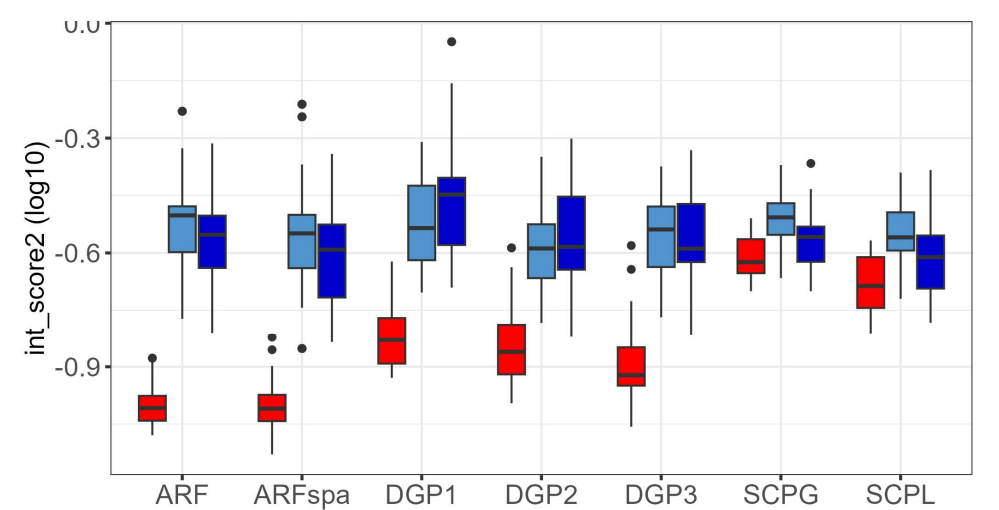
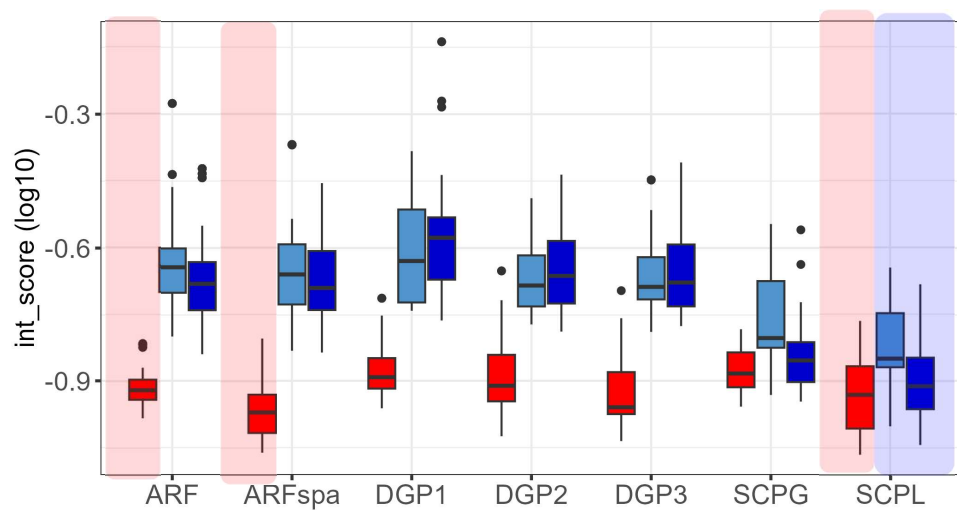
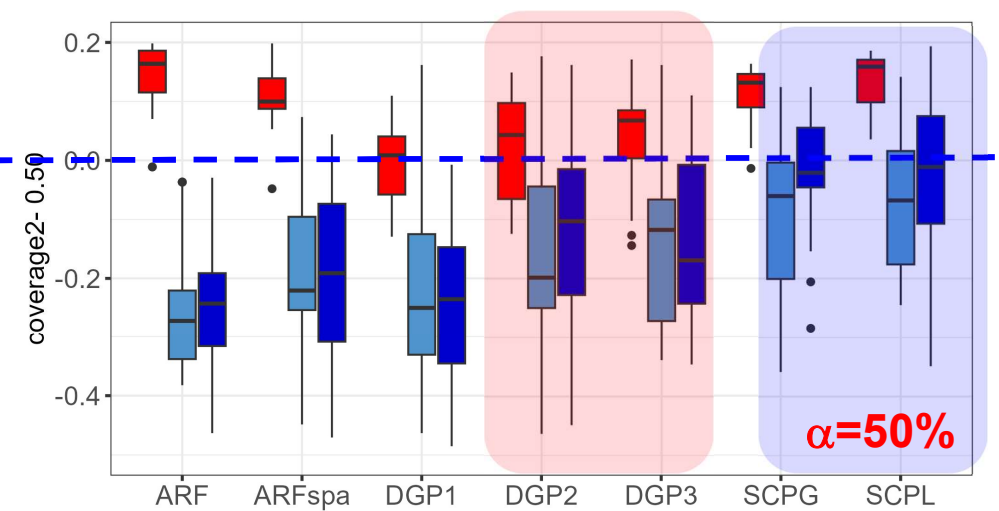
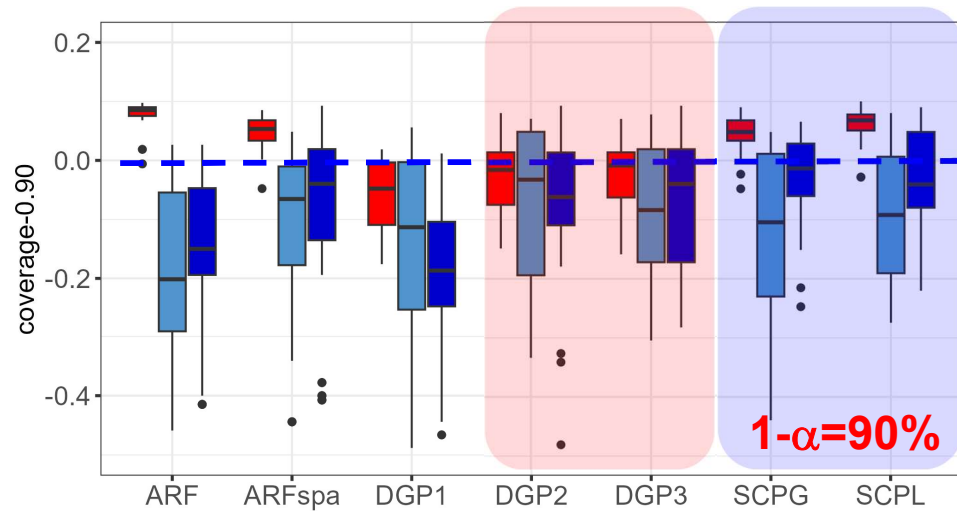
# Robustness to the characteristics of the sample distribution



## Results of 25 repeated random experiments – real case – **N=500**

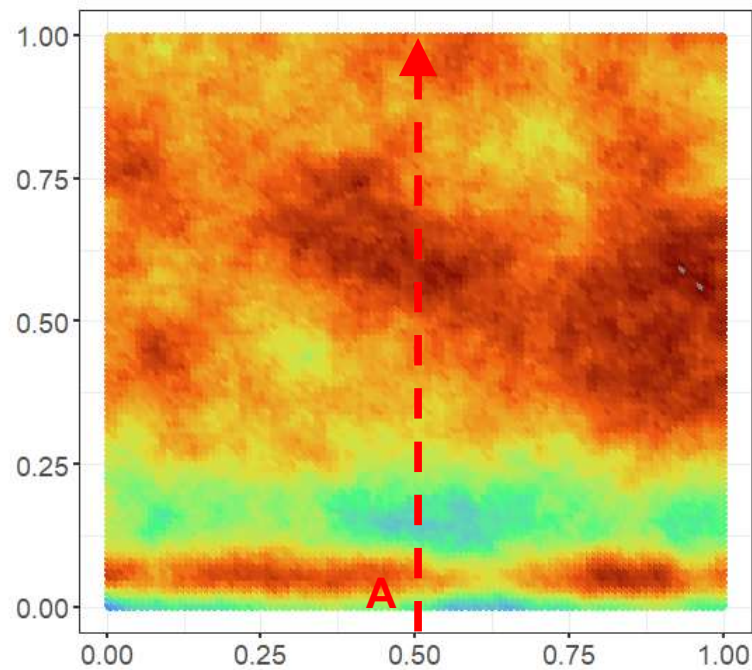


## Results of 25 repeated random experiments – real case – N=125



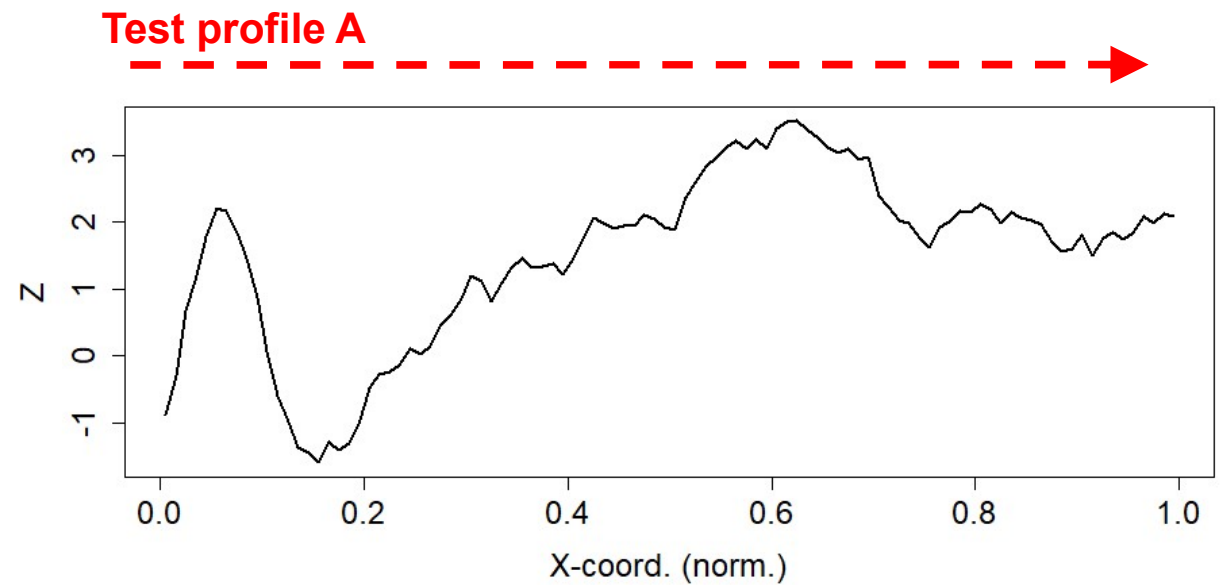
?

## Benchmark synthetic case



Zero-centered 2D Gaussian process with  
spherical covariance (range=0.35,  $\sigma=0.5$ )

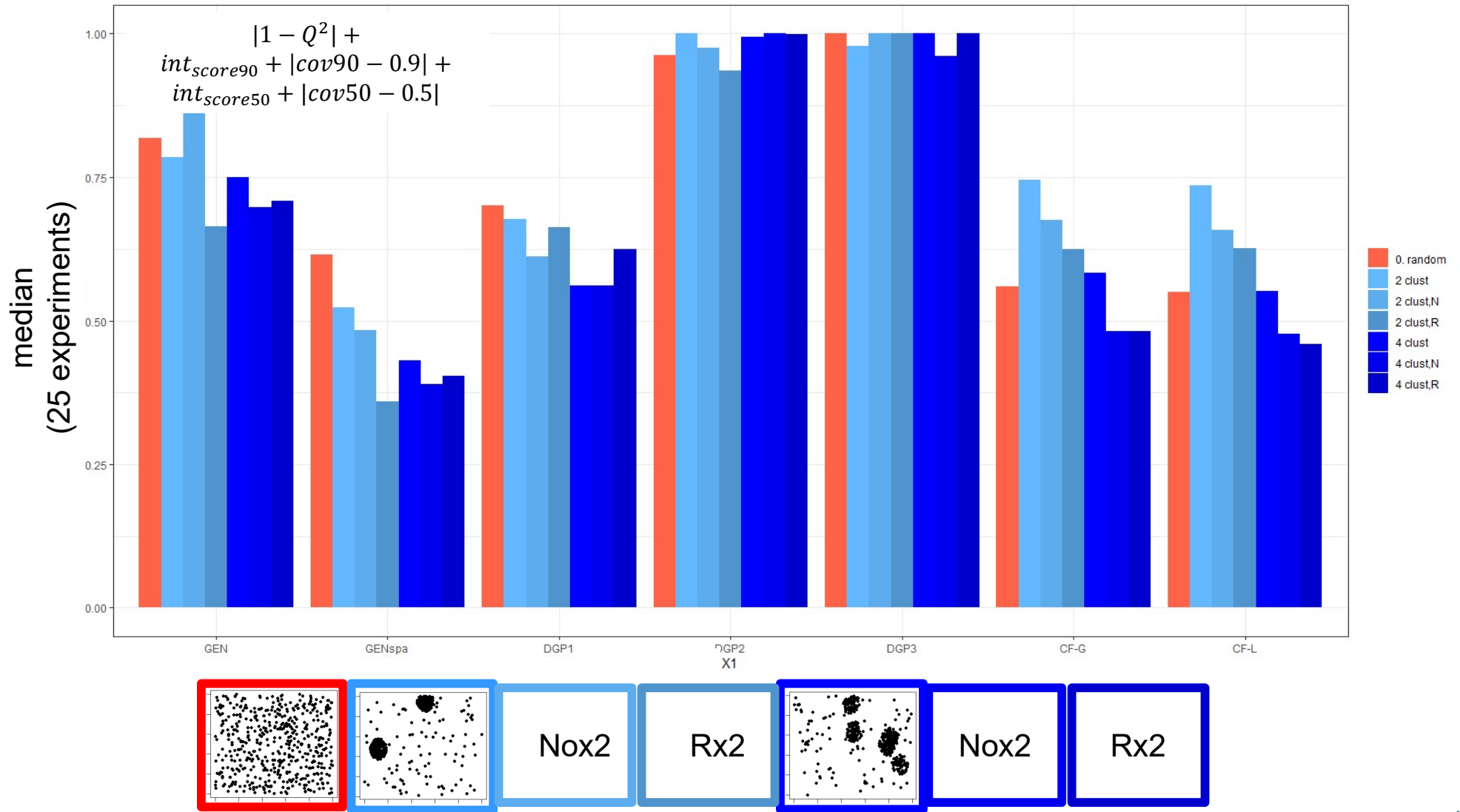
+  
 $X \cdot \sin(X)$



> 48  
> 48



# Robustness to the characteristics of the samples' distribution - **synthetic**



# CV applied to dune case

