

# Nonlinear dimensionality reduction for functional computer code modelling

Benjamin Auder

CEA - UPMC

24 august 2010

Thesis since 02/2008

PhD supervisors : Gérard Biau (UPMC)  
Bertrand Iooss (EDF)

## Industrial context

Framework : life span of reactor vessels.

→ Several sequences of accidents can occur.

*Goal* = estimate their probabilities of occurrence.

# Industrial context

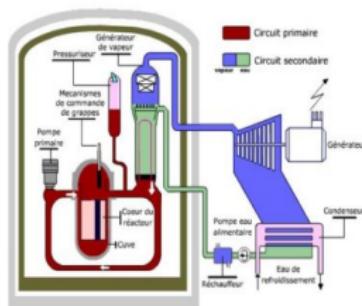
Framework : life span of reactor vessels.

→ Several sequences of accidents can occur.

*Goal* = estimate their probabilities of occurrence.

## Methodology

### Modelling



# Industrial context

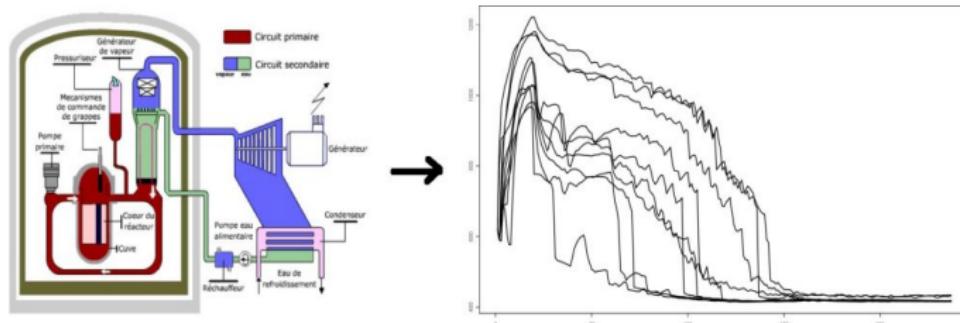
Framework : life span of reactor vessels.

→ Several sequences of accidents can occur.

Goal = estimate their probabilities of occurrence.

## Methodology

Modelling → Simulation



# Industrial context

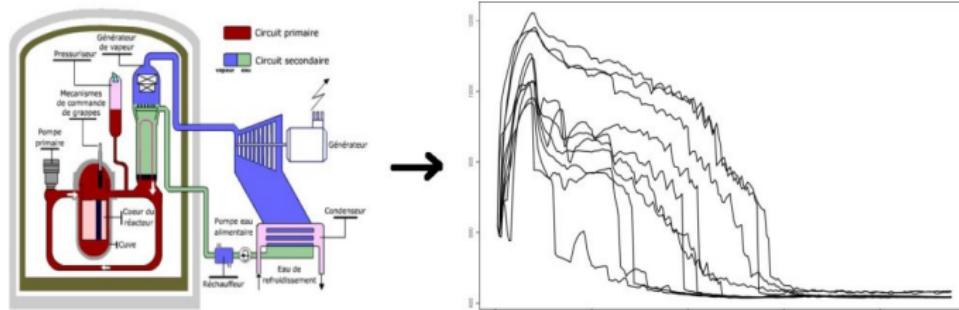
Framework : life span of reactor vessels.

→ Several sequences of accidents can occur.

*Goal* = estimate their probabilities of occurrence.

## Methodology

Modelling → Simulation → Computation.



→ Sensitivity analysis,  
uncertainty propagation  
...

# Industrial context

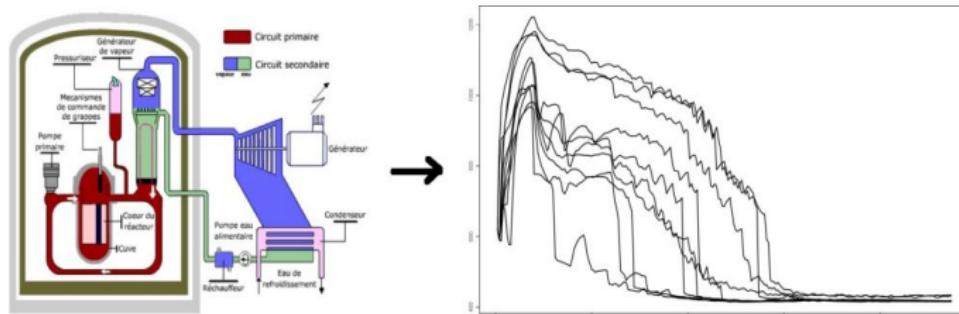
Framework : life span of reactor vessels.

→ Several sequences of accidents can occur.

Goal = estimate their probabilities of occurrence.

## Methodology

Modelling → **Simulation** → Computation.



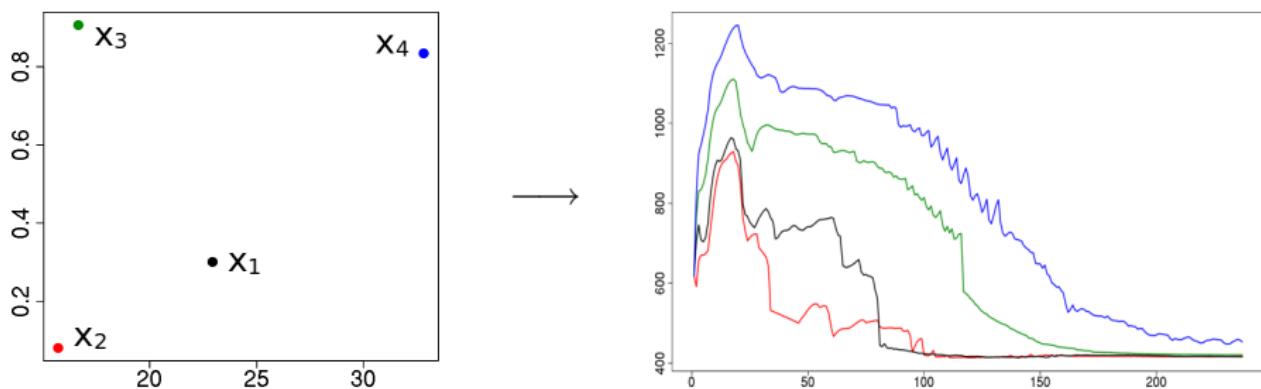
→ Sensitivity analysis,  
uncertainty propagation  
...

*Improve the simulation stage to allow accurate computations*

# Mathematical formulation

$n$  known couples  $(x_i, y_i)$  = inputs-outputs of a **very slow code** :

- Inputs  $x_i \in \mathbb{R}^p$  = initial state of physical system ;
- Outputs  $y_i \in \mathcal{C}([a, b], \mathbb{R})$  = evolutions of parameters.



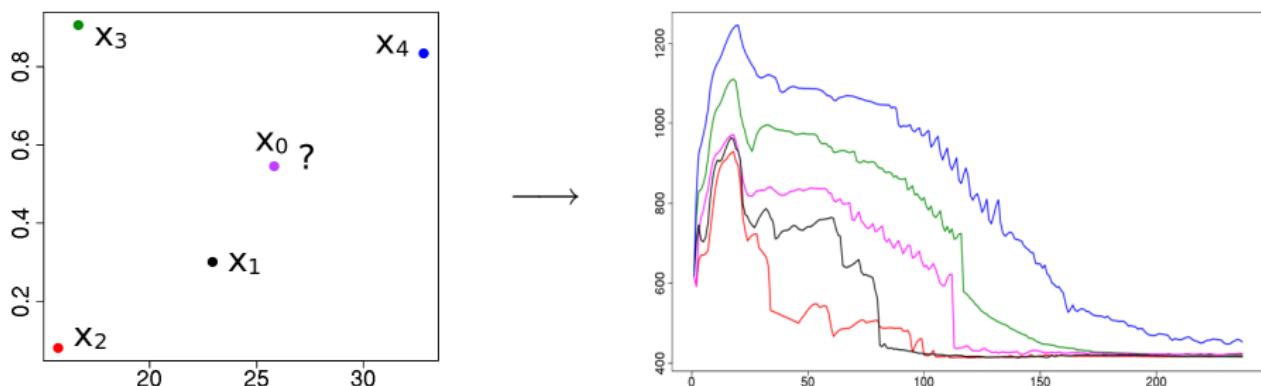
# Mathematical formulation

$n$  known couples  $(x_i, y_i)$  = inputs-outputs of a **very slow code** :

- Inputs  $x_i \in \mathbb{R}^p$  = initial state of physical system ;
- Outputs  $y_i \in \mathcal{C}([a, b], \mathbb{R})$  = evolutions of parameters.

*Goal* = prediction of functional data :

$$y^{\text{new}} \simeq \varphi(x^{\text{new}}).$$



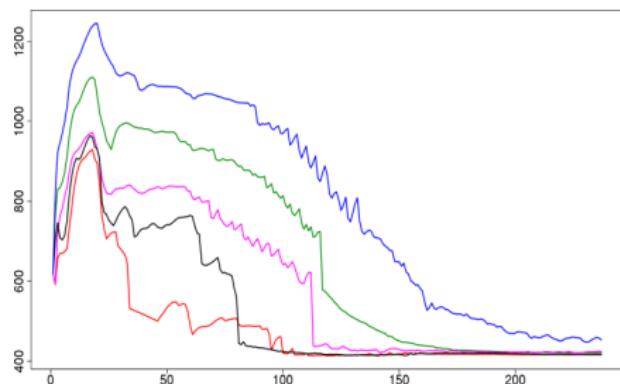
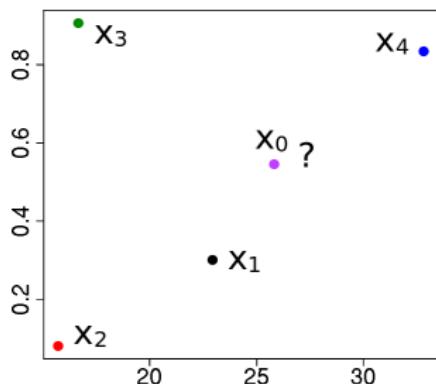
# Mathematical formulation

$n$  known couples  $(x_i, y_i)$  = inputs-outputs of a **very slow code** :

- Inputs  $x_i \in \mathbb{R}^p$  = initial state of physical system ;
- Outputs  $y_i \in \mathcal{C}([a, b], \mathbb{R})$  = evolutions of parameters.

Goal = **prediction** of functional data :

$$y^{\text{new}} \simeq \varphi(x^{\text{new}}).$$



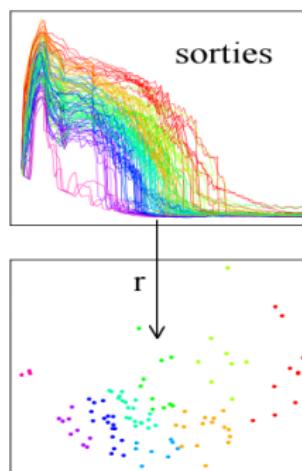
*Statistical learning "regression"  $\mathbb{R}^p \rightarrow \mathcal{C}([a, b], \mathbb{R})$*

Back to the "simple" case of  $y_i \in \mathbb{R}^d$

# Back to the "simple" case of $y_i \in \mathbb{R}^d$

① dimensionality reduction :

$$r : \mathcal{C}([a, b], \mathbb{R}) \rightarrow \mathbb{R}^d \text{ (representation)};$$



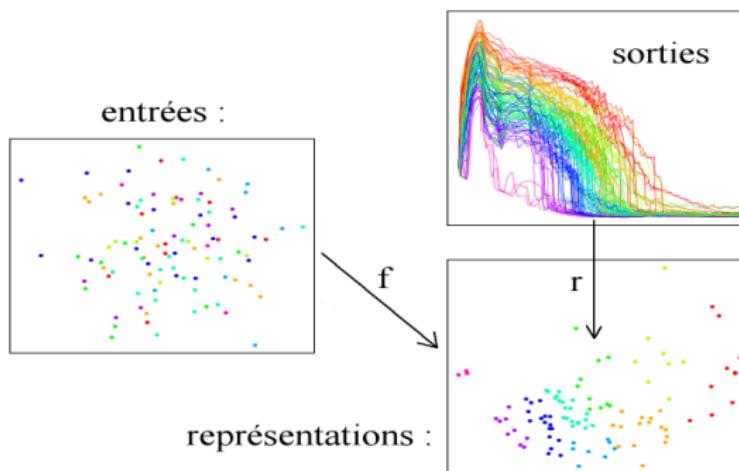
## Back to the "simple" case of $y_i \in \mathbb{R}^d$

- ① dimensionality reduction :

$r : \mathcal{C}([a, b], \mathbb{R}) \rightarrow \mathbb{R}^d$  (representation) ;

- ② statistical learning :

$f : \mathbb{R}^p$  (inputs)  $\rightarrow \mathbb{R}^d$  (reduced outputs) ;



## Back to the "simple" case of $y_i \in \mathbb{R}^d$

- ① dimensionality reduction :

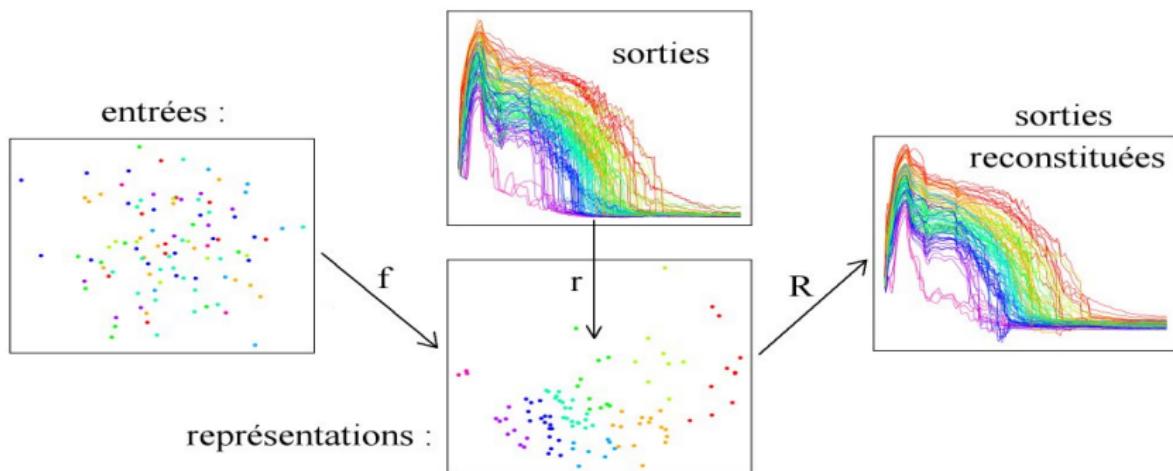
$r : \mathcal{C}([a, b], \mathbb{R}) \rightarrow \mathbb{R}^d$  (representation) ;

- ② statistical learning :

$f : \mathbb{R}^p$  (inputs)  $\rightarrow \mathbb{R}^d$  (reduced outputs) ;

- ③ output space parametrization :

$R : \mathbb{R}^d \rightarrow \mathcal{C}([a, b], \mathbb{R})$  (reconstruction).



# Back to the "simple" case of $y_i \in \mathbb{R}^d$

- ① dimensionality reduction :

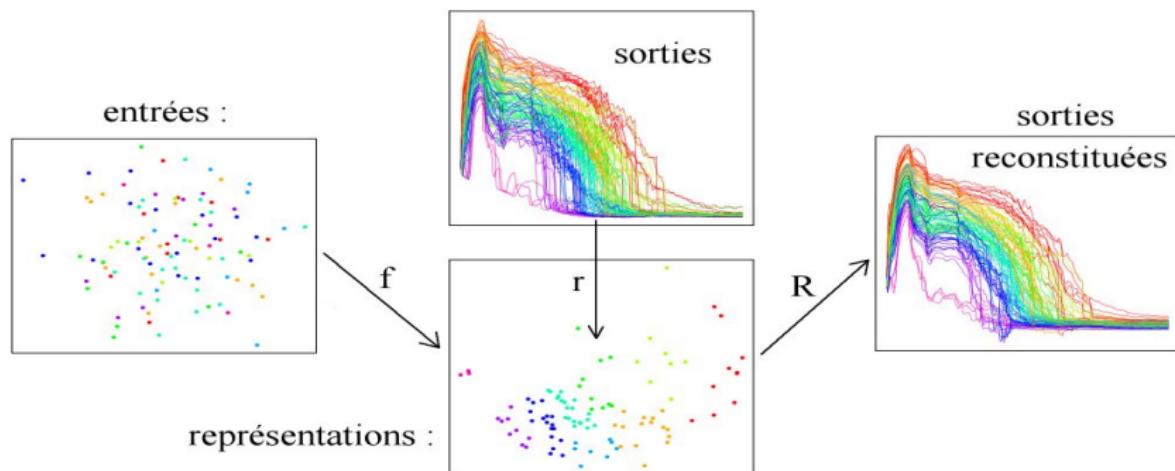
$r : \mathcal{C}([a, b], \mathbb{R}) \rightarrow \mathbb{R}^d$  (representation) ;

- ② statistical learning :

$f : \mathbb{R}^p$  (inputs)  $\rightarrow \mathbb{R}^d$  (reduced outputs) ;

- ③ output space parametrization :

$R : \mathbb{R}^d \rightarrow \mathcal{C}([a, b], \mathbb{R})$  (reconstruction).



# State of art

## (More or less) classical methods

- Functional linear regression :  
Faraway, 1997 ; Ramsay & Silverman, 2005, ...

# State of art

## (More or less) classical methods

- Functional linear regression :  
Faraway, 1997 ; Ramsay & Silverman, 2005, ...
- Decomposition on an orthonormal basis, then learning of  $d$ -dimensional coefficients :  
Chiou et al., 2004 ; Govaerts & Noël, 2005 ;  
Bayarri et al., 2007 ; Marrel, 2008 ; Monestiez & Nerini, 2009

# State of art

## (More or less) classical methods

- Functional linear regression :  
Faraway, 1997 ; Ramsay & Silverman, 2005, ...
- Decomposition on an orthonormal basis, then learning of  $d$ -dimensional coefficients :  
Chiou et al., 2004 ; Govaerts & Noël, 2005 ;  
Bayarri et al., 2007 ; Marrel, 2008 ; Monestiez & Nerini, 2009

(New) goal : minimize the representation dimension  $d$ , to

- simplify the model ;
- avoid overfitting,

while keeping good performances.

## 1 Riemannian Manifold Learning

## 2 Applications

## 1 Riemannian Manifold Learning

## 2 Applications

# Local steps

RML  $\simeq$  preservation of angles *and* geodesic distances.

- ➊ choose an origin curve  $y_0$  among the  $y_i$ , (e.g., the mean);

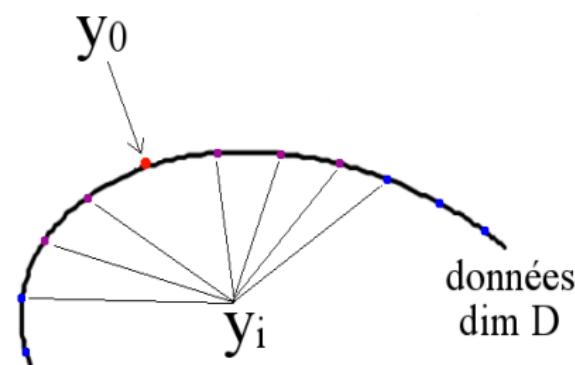


FIG.: Origin curve  $y_0$

# Local steps

RML  $\simeq$  preservation of angles *and* geodesic distances.

- ❶ choose an origin curve  $y_0$  among the  $y_i$ , (e.g., the mean) ;
- ❷ determine a local basis  $Q_0 = (e_1, \dots, e_d)$  for the tangent space at  $y_0$  (PCA on the neighborhoods curves) ;

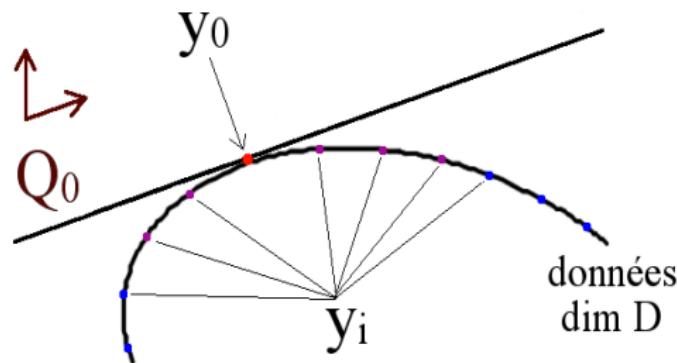


FIG.: Tangent plane at  $y_0$  + local basis  $Q_0$

# Local steps

RML  $\simeq$  preservation of angles *and* geodesic distances.

- ① choose an origin curve  $y_0$  among the  $y_i$ , (e.g., the mean) ;
- ② determine a local basis  $Q_0 = (e_1, \dots, e_d)$  for the tangent space at  $y_0$  (PCA on the neighborhoods curves) ;
- ③ compute the reduced coordinates  $z_i$  or curves  $y_i$  "close" to  $y_0$  by projecting on  $Q_0$ ,

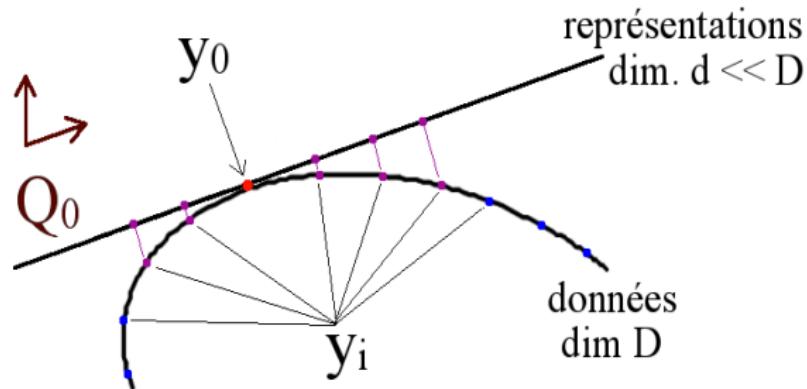


FIG.: Local coordinates  $z_i$  on the tangent plane

# Local steps

RML  $\simeq$  preservation of angles *and* geodesic distances.

- ① choose an origin curve  $y_0$  among the  $y_i$ , (e.g., the mean) ;
- ② determine a local basis  $Q_0 = (e_1, \dots, e_d)$  for the tangent space at  $y_0$  (PCA on the neighborhoods curves) ;
- ③ compute the reduced coordinates  $z_i$  or curves  $y_i$  "close" to  $y_0$  by projecting on  $Q_0$ ,
- ④ then normalize to verify the identity  $\|y - y_0\| = \|x - x_0\|$ .

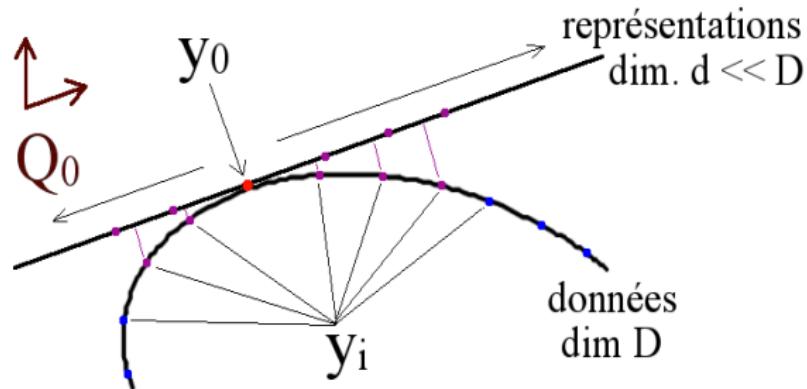
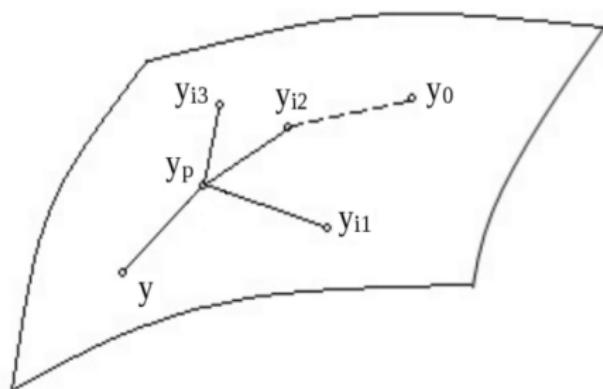


FIG.: Normalization of coordinates  $z_i$

## "far" from $y_0$

Step 4 : for  $y$  far from  $y_0$  (too far for last step to be accurate),

- $y_p$  = predecessor of  $y$  on a shortest path from  $y_0$
- $y_{i_1}, \dots, y_{i_d}$  = neighbors of  $y_p$  for which the  $z_i$  coordinates are known (breadth-first)



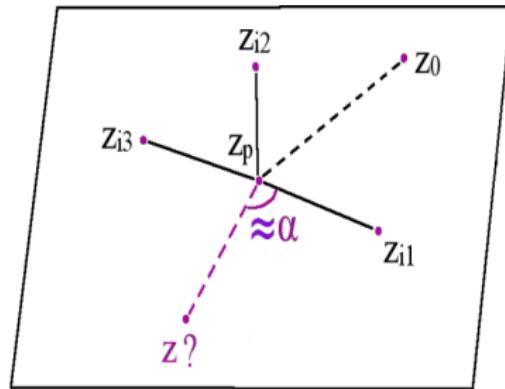
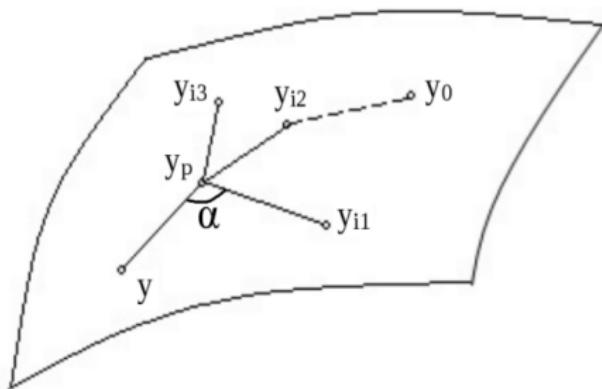
## "far" from $y_0$

Step 4 : for  $y$  far from  $y_0$  (too far for last step to be accurate),

- $y_p$  = predecessor of  $y$  on a shortest path from  $y_0$
- $y_{i_1}, \dots, y_{i_d} =$  neighbors of  $y_p$  for which the  $z_i$  coordinates are known (breadth-first)

$z = r(y)$  computed by..

- ..preserving angles as much as possible :  $\widehat{zz_pz_{i_j}} \simeq \widehat{yy_py_{i_j}}$  ;
- ..under the normalization constraint  $\|y - y_p\| = \|z - z_p\|$ .



# Examples

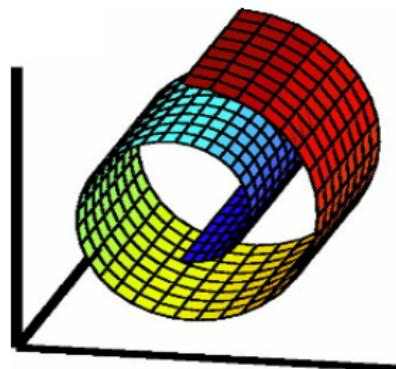


FIG.: 3D Swissroll, 400 points

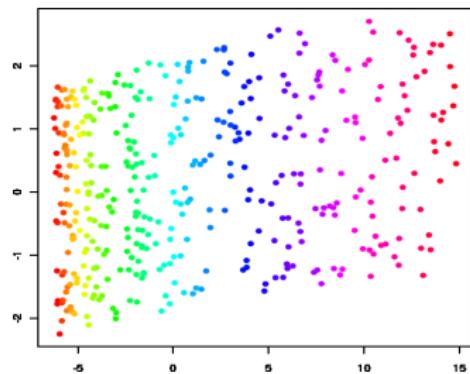


FIG.: RML Representation

# Examples

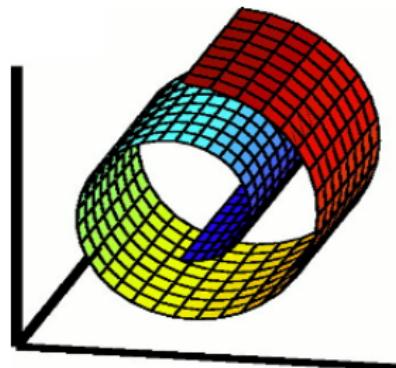


FIG.: 3D Swissroll, 400 points

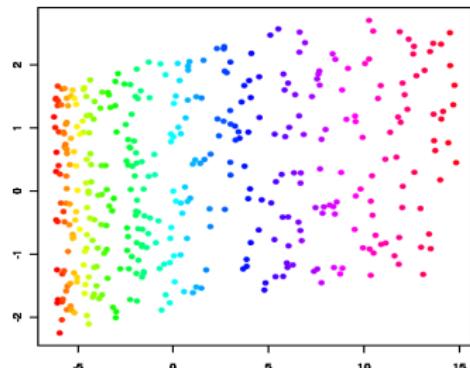


FIG.: RML Representation

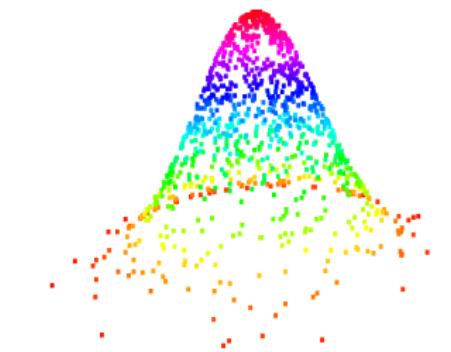


FIG.: 3D Gaussian, 1000 points

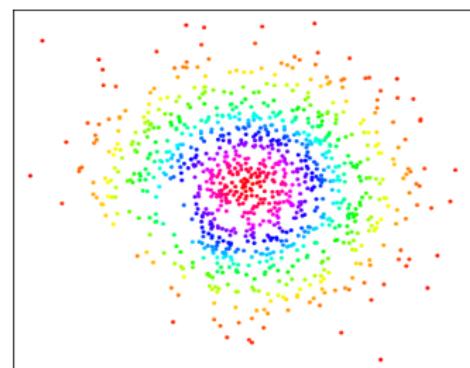


FIG.: RML Representation

## 1 Riemannian Manifold Learning

## 2 Applications

# Validation step

Data :

- training =  $\{(x_i, y_i), i = 1, \dots, n\}$ ;
- test =  $\{(x'_i, y'_i), i = 1, \dots, m\}$ ;

Model predictions :  $\hat{y}'_i = M(x'_i)$ ,  $i = 1, \dots, m$ .

# Validation step

Data :

- training =  $\{(x_i, y_i), i = 1, \dots, n\}$ ;
- test =  $\{(x'_i, y'_i), i = 1, \dots, m\}$ ;

Model predictions :  $\hat{y}'_i = M(x'_i)$ ,  $i = 1, \dots, m$ .

"Absolute" then relative measures of the pointwise error

$$MSE[j] = \frac{1}{m} \sum_{i=1}^m (\hat{y}'_i(j) - y'_i(j))^2, \quad j = 1, \dots, D \text{ (discretization).}$$

# Validation step

Data :

- training =  $\{(x_i, y_i), i = 1, \dots, n\}$ ;
- test =  $\{(x'_i, y'_i), i = 1, \dots, m\}$ ;

Model predictions :  $\hat{y}'_i = M(x'_i)$ ,  $i = 1, \dots, m$ .

"Absolute" then relative measures of the pointwise error

$$MSE[j] = \frac{1}{m} \sum_{i=1}^m (\hat{y}'_i(j) - y'_i(j))^2, \quad j = 1, \dots, D \text{ (discretization).}$$

$$Q_2[j] = 1 - \frac{m.MSE[j]}{\sum_{i=1}^m (\bar{y}(j) - y'_i(j))^2} \text{ (comparison to the mean).}$$

$-\infty < Q_2 \leq 1$  :  
 $\leq 0 \Rightarrow \text{(very) bad model};$   
 $\simeq 1 \Rightarrow \text{perfect model.}$

# Test I - temperature curves

100 model runs,  
4 dimensions on input,  
168 discretization points.

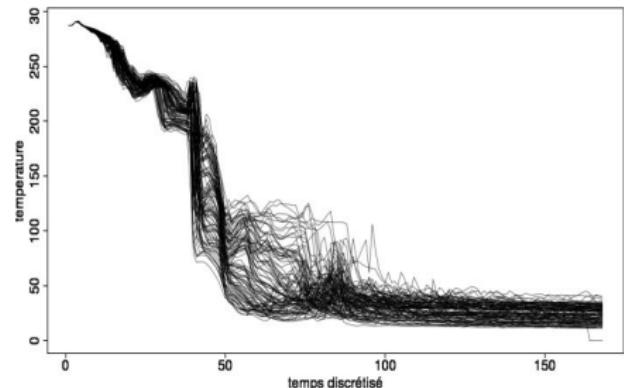


FIG.: The 100 code outputs

# Test I - temperature curves

100 model runs,  
4 dimensions on input,  
168 discretization points.

cross validation  
leave-10-out :

MSE at l.,  $Q_2$  at r.;  $d = 4$

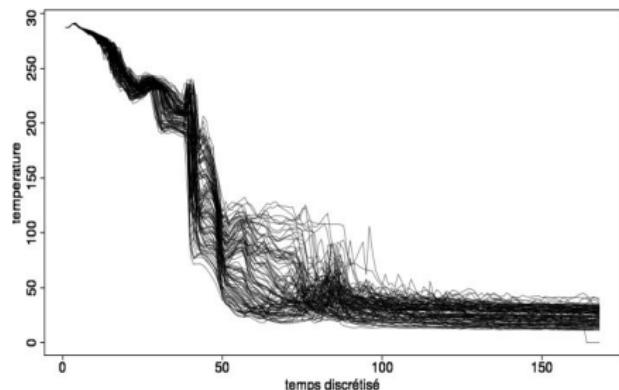


FIG.: The 100 code outputs

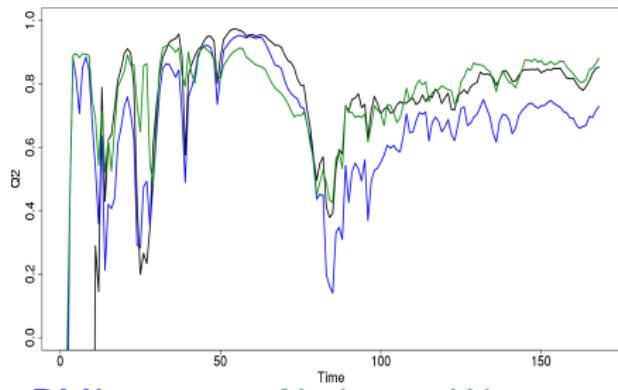
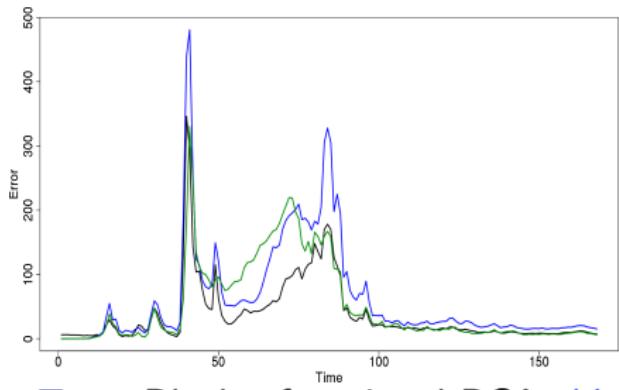


FIG.: Black : functional PCA ; blue : RML ; green : Nadaraya-Watson.

# 5 predicted curves

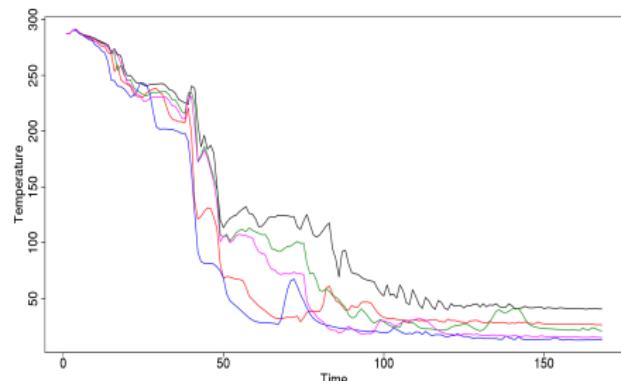


FIG.: Real curves

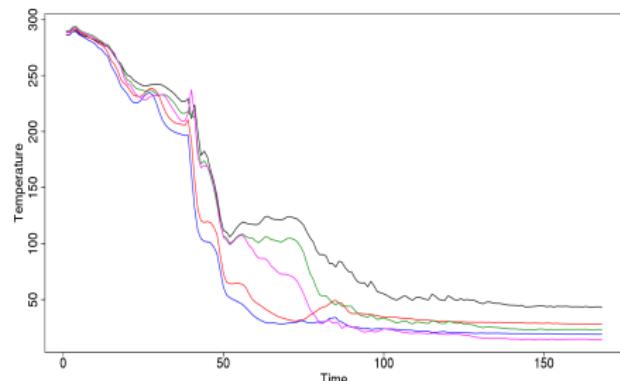


FIG.: PCA dim. red.

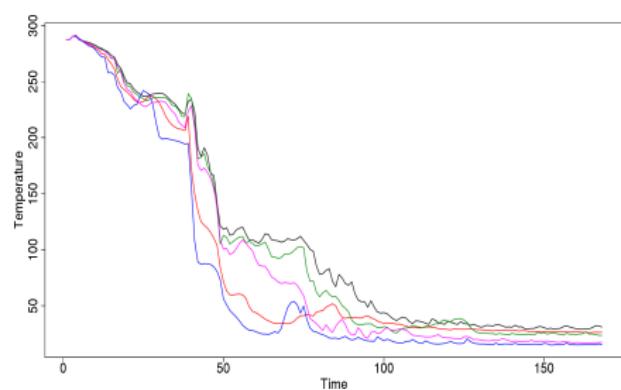


FIG.: RML dim. red.

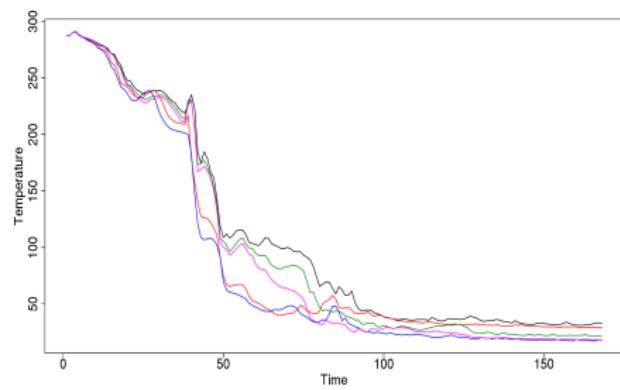


FIG.: Nadaraya-Watson

## Test II - temperature curves

600 model runs,  
11 dimensions on input,  
414 discretization points.

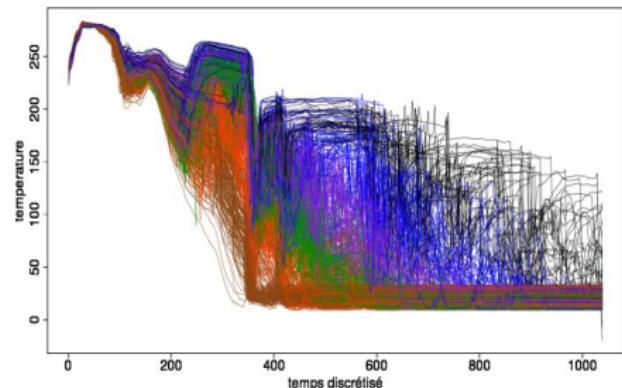


FIG.: The 600 code outputs

## Test II - temperature curves

600 model runs,  
11 dimensions on input,  
414 discretization points.

cross validation  
leave-10-out :

MSE at l.,  $Q_2$  at r.;  $d = 7$

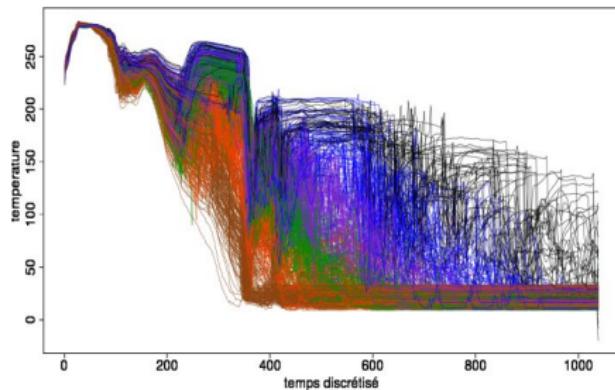


FIG.: The 600 code outputs

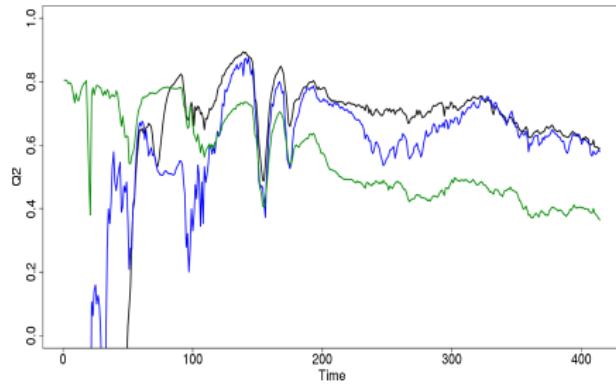
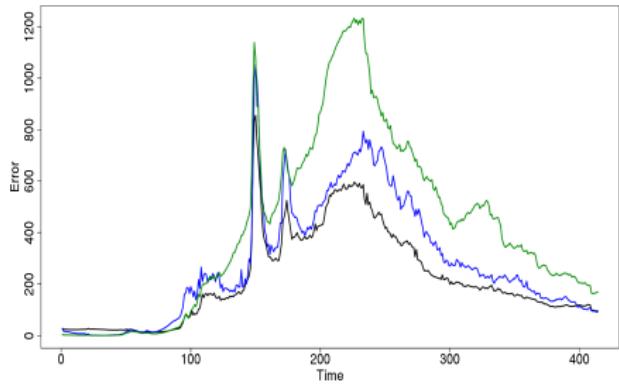


FIG.: Black : functional PCA ; blue : RML ; green : Nadaraya-Watson.

# 5 predicted curves

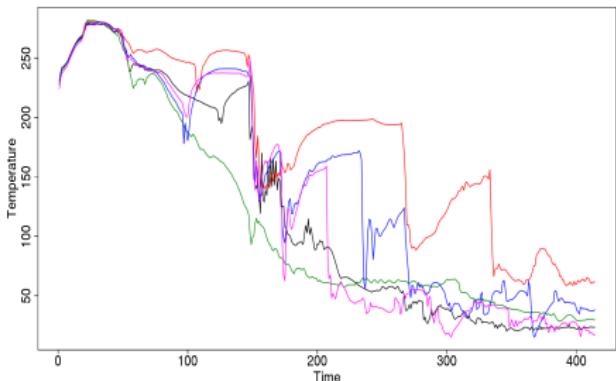


FIG.: Real curves

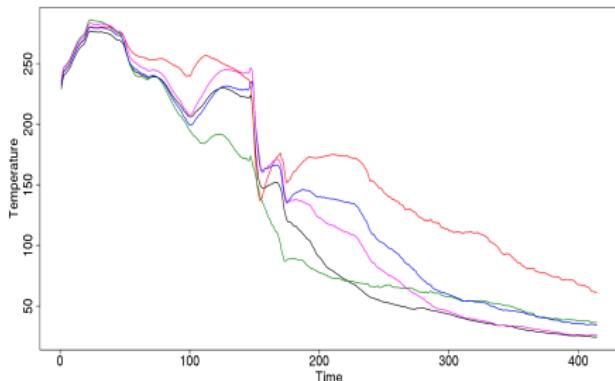


FIG.: PCA red. dim.

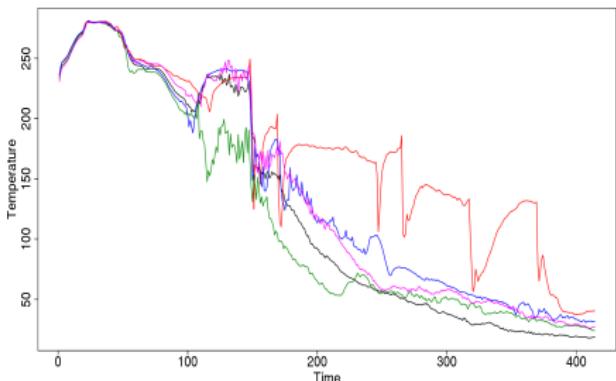


FIG.: RML red. dim.

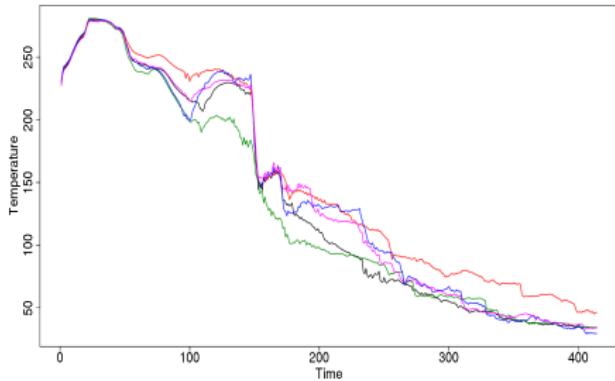


FIG.: Nadaraya-Watson

# Conclusion

*The model is good enough regarding to the industrial context.*

⇒ help to the projet DDVCV (life span of reactor vessels).

# Conclusion

*The model is good enough regarding to the industrial context.*

⇒ help to the projet DDVCV (life span of reactor vessels).

## Dimensionality reduction

Linear (PCA) and non-linear : complementary.

# Conclusion

*The model is good enough regarding to the industrial context.*

⇒ help to the projet DDVCV (life span of reactor vessels).

## Dimensionality reduction

Linear (PCA) and non-linear : complementary.

Missing : proof of convergence to the right manifold . . .

. . . and some parameters which should be "optimized" automatically.

# Conclusion

*The model is good enough regarding to the industrial context.*

⇒ help to the projet DDVCV (life span of reactor vessels).

## Dimensionality reduction

Linear (PCA) and non-linear : complementary.

Missing : proof of convergence to the right manifold . . .

. . . and some parameters which should be "optimized" automatically.

Future research : "functional" principal curves and surfaces.

Example of a principal  
surface in 2D :

