

Modeling Uncertainties by Simulation

Nabil Rachdi (PhD Student)

nabil.rachdi@eads.net

► GdR MASCOT-NUM (17, 18, 19 of March – Avignon) ◀

Advisors:

*Jean-Claude Fort (Paris V), Thierry Klein (Toulouse III)
Fabien Mangeant (EADS IW), Régis Lebrun (EADS IW).*



Outline

1 Motivations and notations

2 Model Selection

3 Risk excess bound

4 Examples

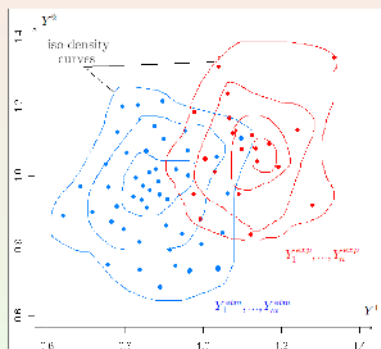
5 Future Work

General framework

- **Variable of Interest:** Y , density f , probability measure \mathbb{Q} **unknown**
- **Quantity of Interest (QoI):** Density distribution, mean, threshold probability, quantile etc...
- **Experimental data:** $\mathcal{Y}_n^{exp} = Y_1^{exp}, \dots, Y_n^{exp}$ (*a priori training data*) supposed i.i.d from \mathbb{Q}
 - Link to history
 - Arise from experiments, complex codes etc...
 - Small number
 - Difficult to obtain
- **Simulated data:** $\mathcal{Y}_m^{sim} = Y_1^{sim}, \dots, Y_m^{sim}$ depend on the model h and the parameter θ
 - $Y_i^{sim} = h(X_i, \theta)$, $i = 1, \dots, m$, X_i i.i.d random variables (density p_X).
 - $h \in \mathcal{H}$ (set of models), $\theta \in \Theta$ (set of parameters)
- **Goal:** Use Simulated data to improve QoI estimation of Y :
 - ⇒ 1. Calibration procedure: choice of the model h , and parameter θ
 - ⇒ 2. Study of the QoI based on $(Y_1^{exp}, \dots, Y_n^{exp}, \hat{Y}_1^{sim}, \dots, \hat{Y}_m^{sim})$ (*a posteriori training data*) to compare with QoI based on $(Y_1^{exp}, \dots, Y_n^{exp})$

Simulated data calibration

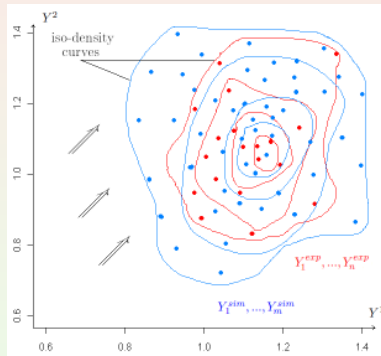
- Illustration of Experimental & Simulated Data: Example of a 2D-Performance $Y = (Y^1, Y^2)$



- Choice of $h \in \mathcal{H}$ and $\theta \in \Theta$? \Rightarrow driven by the QoI

Simulated data calibration

- Illustration of Experimental & Simulated Data: Example of a 2D-Performance $Y = (Y^1, Y^2)$



- Choice of $h \in \mathcal{H}$ and $\theta \in \Theta$? \Rightarrow driven by the QoI

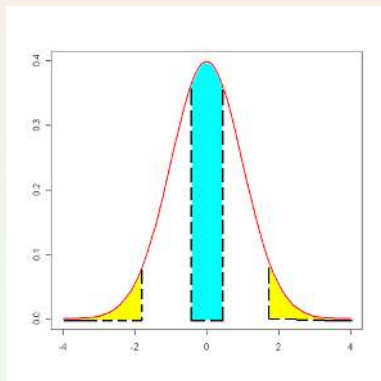
Quantity of Interest (QoI):

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space, \mathbb{D} a metric space, and W a random variable defined on $(\Omega, \mathcal{A}, \mathbb{P})$, the QoI of W is defined as the function

$$\begin{aligned} \varphi : (\Omega, \mathcal{A}, \mathbb{P}) &\longrightarrow \mathbb{D} \\ W &\longmapsto \varphi(W). \end{aligned}$$

Some QoI:

- $\varphi(W) = \mathbb{E}(W), q_W^\alpha \Rightarrow \mathbb{D} = \mathbb{R}$
- $\varphi(W) = \mathbb{P}(W > s) \Rightarrow \mathbb{D} = [0, 1]$
- $\varphi(W) = f_W \Rightarrow \mathbb{D} = \{\text{set of distributions}\}$



Choice of (h, θ) for a QoI φ

- **Minimization of a criterion :**

$$M(h, \theta) = \mathcal{D}(\varphi_{h, \theta}, \varphi_Y), \quad \mathcal{D} : \text{distance on } \mathbb{D} \times \mathbb{D}$$

- $\varphi_{h, \theta}$ and φ_Y are QoI of $h(X, \theta)$ and Y (resp.)
- suppose (h^*, θ^*) is the **unique** minimum of M

- **GOAL:**

- (Current work) Minimize $M(h, \theta)$ over Θ for **fixed** $h \in \mathcal{H}$

$$\theta_0(h) = \underset{\theta \in \Theta}{\text{Argmin}} M(h, \theta)$$

- (Later) Minimize $M(h, \theta_0(h))$ over \mathcal{H}

● We use the form:

$$M(h, \theta) = \mathcal{D}(\varphi_{h,\theta}, \varphi_Y) \longrightarrow M(h, \theta) = \int_{\mathbb{R}} \gamma_{h,\theta}(y) f(y) dy$$

- the function $\gamma_{h,\theta}$ is called **contrast** of (h, θ) :
 $\gamma_{h,\theta} = \Psi(\varphi_{h,\theta})$ with Ψ some function
- recall that f is the density of Y (*unknown*)

● Example of contrasts:

If the *QoI* is the density $\varphi_{h,\theta} = f_{h,\theta}$

- $y \mapsto \gamma_{h,\theta}(y) = -\ln(f_{h,\theta}(y)) \Rightarrow M(h, \theta) = K(f, f_{h,\theta})$
- $y \mapsto \gamma_{h,\theta}(y) = \|f_{h,\theta}\|_2^2 - 2f_{h,\theta}(y) \Rightarrow M(h, \theta) = \|f - f_{h,\theta}\|_2^2$.
- etc...

If the *QoI* is the mean $\varphi_{h,\theta} = \mathbb{E}_X(h(X, \theta))$

- $y \mapsto \gamma_{h,\theta}(y) = (y - \mathbb{E}_X(h(X, \theta)))^2 \Rightarrow M(h, \theta) = \mathbb{E}_Y(Y - \mathbb{E}_X(h(X, \theta)))^2$

Etc ...

- **Criterion to minimize:**

$$M(h, \theta) = \int_{\mathbb{R}} \gamma_{h, \theta}(y) f(y) dy$$

- **Difficulties:**

- The density function f of Y is **unknown**
- For complex models h , the QoI $\varphi_{h, \theta}$ can be **unreachable** with reasonable CPU time $\rightarrow \gamma_{h, \theta} = \Psi(\varphi_{h, \theta})$ **unreachable**.

- **Alternative: Use of Experimental & Simulated data**

- Replace f by its *empirical* version $\rightarrow \frac{1}{n} \sum_{i=1}^n \delta_{Y_i^{exp}}$ (depends on n -Experimental data \mathcal{Y}_n^{exp})
- Replace $\gamma_{h, \theta}$ (precisely $\varphi_{h, \theta}$) by its *simulated* version $\rightarrow \gamma_{h, \theta}^m = \Psi(\varphi_{h, \theta}^m)$ (depends on m -Simulated data \mathcal{Y}_m^{sim}).

- **Practical criterion**

$$M(h, \theta) = \int_{\mathbb{R}} \gamma_{h, \theta}(y) f(y) dy \quad \longleftrightarrow \quad M_{n, m}(h, \theta) := \frac{1}{n} \sum_{i=1}^n \gamma_{h, \theta}^m(Y_i^{exp})$$

- **Criterion to minimize in practice:**

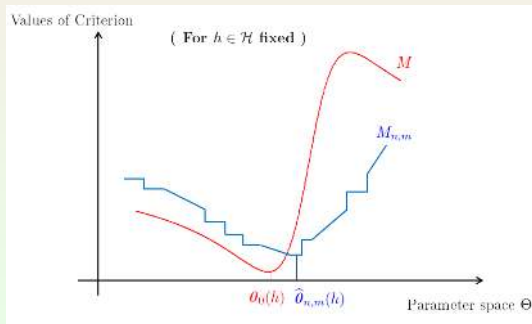
$$M_{n,m}(h, \theta) := \frac{1}{n} \sum_{i=1}^n \gamma_{h,\theta}^m(Y_i^{\text{exp}})$$

- **Estimator of $\theta_0(h) = \text{Argmin}_{\theta \in \Theta} M(h, \theta)$**

$$\hat{\theta}_{n,m}(h) = \text{Argmin}_{\theta \in \Theta} M_{n,m}(h, \theta).$$

First question: **Consistency**

$$\hat{\theta}_{n,m}(h) \xrightarrow[n \rightarrow +\infty]{m \rightarrow +\infty} \theta_0(h)?$$



Proposition: Oracle inequality

We prove

$$\underbrace{M(h, \hat{\theta}_{n,m}(h)) - M(h^*, \theta^*)}_{\text{risk excess of } (h, \hat{\theta}_{n,m}(h))} \leq 2 \cdot \frac{1}{\sqrt{n}} \|\mathbb{G}_n \gamma_{h,\cdot}^m\|_{\Theta} + 2 \cdot \|\mathcal{E}_h^m\|_{\Theta} + \Delta_h$$

• Variance terms :

$$- \frac{1}{\sqrt{n}} \|\mathbb{G}_n \gamma_{h,\cdot}^m\|_{\Theta} = \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \left(\gamma_{h,\theta}^m(Y_i^{\text{exp}}) - \mathbb{E}_Y(\gamma_{h,\theta}^m(Y)) \right) \right| \quad (\text{deviation})$$

\Rightarrow *Estimation Error of Statistical Data* \rightarrow depends on contrast (i.e QoI)

$$- \|\mathcal{E}_h^m\|_{\Theta} = \sup_{\theta \in \Theta} \|\gamma_{h,\theta}^m - \gamma_{h,\theta}\|_{1, \mathbb{Q}}, \quad \text{with} \quad \|g\|_{1, \mathbb{Q}} = \int_{\mathbb{R}} |g(y)| f(y) dy$$

\Rightarrow *Simulation Error*

• Bias term :

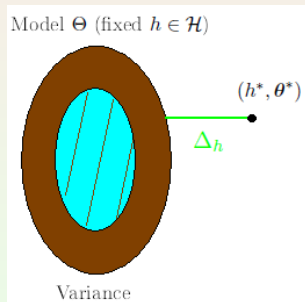
$$- \Delta_h = M(h, \theta_0(h)) - M(h^*, \theta^*)$$

\Rightarrow *Approximation Error of the model h*

Representation of the errors

- **Risk minimization :**

$$M(h, \hat{\theta}_{n,m}) - M(h^*, \theta^*) \leq 2 \cdot \frac{1}{\sqrt{n}} \|\mathbb{G}_n \gamma_{h^*}^m\|_{\Theta} + 2 \cdot \|\mathcal{E}_h^m\|_{\Theta} + \Delta_h$$



- **Work:** study the **simulation effect** on the calibration procedure.

Oracle Inequality

Recall $M(h, \theta) = \int_{\mathbb{R}} \gamma_{h, \theta}(y) f(y) dy$, we have

$$M(h, \hat{\theta}_n(h)) - M(h^*, \theta^*) \leq \frac{1}{\sqrt{n}} \|\mathbb{G}_n \gamma_{h, \cdot}\|_{\Theta} + \Delta_h$$

- **For no complex models:** (Linear model etc...)

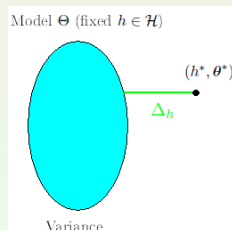
\Rightarrow the QoI $\varphi_{h, \theta}$ is **reachable** \Rightarrow Simulation is **useless**

Advantages

- Well studied
- Maximum Likelihood etc...

Drawback

- Δ_h can be large (due to simplification of h)
Trade-off Bias-Variance



- **For Input/Output data:** $Z_1 = (X_1, Y_1^{exp}), \dots, Z_n = (X_n, Y_n^{exp})$

$z = (x, y) \quad \gamma_{h, \theta}(z) = (y - h(x, \theta))^2 \Rightarrow$ Simulation is **useless**

Theorem: Consistency

We prove that $\hat{\theta}_{n,m}(h) \xrightarrow[n \rightarrow +\infty]{m \rightarrow +\infty} \theta_0(h)$, under some conditions in terms of :

- Model Complexity (*Bracketing Numbers*) \Rightarrow depending on the *quantity of interest*
- Simulation Speed (*Size of Simulated Data set, m*)
- Control of Simulated contrasts (*Modified Lindeberg conditions*)

• **Consequence:** Compute a *QoI* based on $(Y_1^{exp}, \dots, Y_n^{exp}, \hat{Y}_1^{sim}, \dots, \hat{Y}_m^{sim})$ is better than a *QoI* based on $(Y_1^{exp}, \dots, Y_n^{exp})$ in some typical cases.

- Key tool : -*Empirical process theory*- V.d.Vaart (1996,2000), V.d Geer (2000) etc ...
- We give practical conditions for a wide range of applications.

Theorem

Let $\Gamma_h^m := \{\gamma_{h,\theta}^m, \theta \in \Theta\}$ and denote by F_m an envelope function, assume that

- $R_h^m(\theta, \theta') = \mathbb{Q} \gamma_{h,\theta}^m \gamma_{h,\theta'}^m - \mathbb{Q} \gamma_{h,\theta}^m \mathbb{Q} \gamma_{h,\theta'}^m$ converges on $\Theta \times \Theta$
- $\sup_{d(\theta, \theta') \leq \delta_m} \mathbb{Q} \left(\gamma_{h,\theta}^m - \gamma_{h,\theta'}^m \right)^2 \xrightarrow{m \rightarrow +\infty} 0, \quad \forall \delta_m \downarrow 0$
- (i) $\mathbb{Q} F_m^2 = O(1)$
- (ii) $\mathbb{Q} F_m^2 \mathbb{1} \{F_m > \sqrt{n} \epsilon\} \xrightarrow{n, m \rightarrow +\infty} 0 \quad \forall \epsilon > 0.$

$$J_{[\cdot]}(\delta_m, \Gamma_h^m, L_2(\mathbb{Q})) \xrightarrow{m \rightarrow +\infty} 0, \quad \forall \delta_m \downarrow 0,$$

then $\|\mathbb{G}_n\|_{\Gamma_h^m}$ converges ($n, m \rightarrow +\infty$) to the supremum of a centered Gaussian process with covariance function

$$R_h(\theta, \theta') = \mathbb{Q} \gamma_{h,\theta} \gamma_{h,\theta'} - \mathbb{Q} \gamma_{h,\theta} \mathbb{Q} \gamma_{h,\theta'}.$$

Corollary

If Γ_h^m satisfies to conditions of this Theorem, the calibration procedure is *consistent*, i.e

$$d\left(\widehat{\theta}_{n,m}(h), \theta_0(h)\right) \xrightarrow{n, m \rightarrow +\infty} 0.$$

Theorem

Let $\Gamma_h^m := \{\gamma_{h,\theta}^m, \theta \in \Theta\}$ and denote by F_m an envelope function, assume that

- $R_h^m(\theta, \theta') = \mathbb{Q} \gamma_{h,\theta}^m \gamma_{h,\theta'}^m - \mathbb{Q} \gamma_{h,\theta}^m \mathbb{Q} \gamma_{h,\theta'}^m$ converges on $\Theta \times \Theta$
- $\sup_{d(\theta, \theta') \leq \delta_m} \mathbb{Q} \left(\gamma_{h,\theta}^m - \gamma_{h,\theta'}^m \right)^2 \xrightarrow{m \rightarrow +\infty} 0, \quad \forall \delta_m \downarrow 0$
- (i) $\mathbb{Q} F_m^2 = O(1)$
- (ii) $\mathbb{Q} F_m^2 \mathbb{1} \{F_m > \sqrt{n} \epsilon\} \xrightarrow{n, m \rightarrow +\infty} 0 \quad \forall \epsilon > 0$ (*Control of Simulated contrasts*)
- $J_{[\cdot]}(\delta_m, \Gamma_h^m, L_2(\mathbb{Q})) \xrightarrow{m \rightarrow +\infty} 0, \quad \forall \delta_m \downarrow 0$ (*Complexity condition*)

then $\|\mathbb{G}_n\|_{\Gamma_h^m}$ converges ($n, m \rightarrow +\infty$) to the supremum of a centered Gaussian process with covariance function

$$R_h(\theta, \theta') = \mathbb{Q} \gamma_{h,\theta} \gamma_{h,\theta'} - \mathbb{Q} \gamma_{h,\theta} \mathbb{Q} \gamma_{h,\theta'}.$$

Corollary

If Γ_h^m satisfies to conditions of this Theorem, the calibration procedure is *consistent*, i.e

$$d\left(\widehat{\theta}_{n,m}(h), \theta_0(h)\right) \xrightarrow[n, m \rightarrow +\infty]{\mathbb{P}} 0.$$

- Risk excess \leq Variance terms + Bias term
- What could happen ?
 - On one hand, a numerician only focuses on minimizing the *bias term* (Δ_h),
 - on the other hand, a statistician can control the *variance term* and ignore the *bias term*.
- We propose a **simultaneous approach** driven by Simulations
 - ⇒ Control of variability + Representativity of the model h
- Consequences :
 - The variance ($\Rightarrow \frac{1}{\sqrt{n}} \|\mathbb{G}_n \gamma_{h,\cdot}^m\|_{\Theta} + \|\mathcal{E}_h^m\|_{\Theta}$) depends on
 - ⇒ the Experimental data
 - ⇒ the *Quantity of Interest* (contrast)
 - ⇒ and the Simulated data
 - We expect a better estimation procedure for **limited amount of experimental data** and **complex models**.

Example of the *Range* study

Phenomenon : $Y = \text{Range}$ (distance an aircraft can travel), **QoI** = density distribution

- **A priori training data** : Experimental data, $n = 20$, $\mathcal{Y}_n^{\text{exp}} = Y_1^{\text{exp}}, \dots, Y_n^{\text{exp}}$
(obtained from complex model h^* supposed to be the "true")

- **Additional knowledge** : Simulated data, $m = 3000$, $\mathcal{Y}_m^{\text{sim}} = Y_1^{\text{sim}}, \dots, Y_m^{\text{sim}}$ from

$$h(X, \theta) = \frac{F V}{C_s} \frac{1}{\theta_1} \log \left(\frac{1}{1 - \theta_2} \right)$$

- Uncertain Inputs $X = (F, V, C_s)^T$
- Parameters $\theta = (\theta_1, \theta_2) \in \Theta$

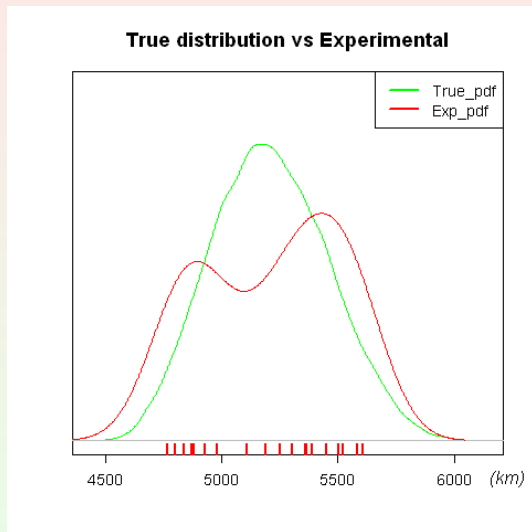
- **Choice of θ ?** $\theta_0(h) = \underset{\theta \in \Theta}{\text{Argmin}} \int_{\mathbb{R}} \gamma_{h, \theta}(y) f(y) dy$

$$\gamma_{h, \theta} = -\ln(f_{h, \theta}) \quad f_{h, \theta} \leftrightarrow f_{h, \theta}^m \text{ (Kernel)} \quad f \leftrightarrow \frac{1}{n} \sum_{i=1}^n \delta_{Y_i^{\text{exp}}}$$

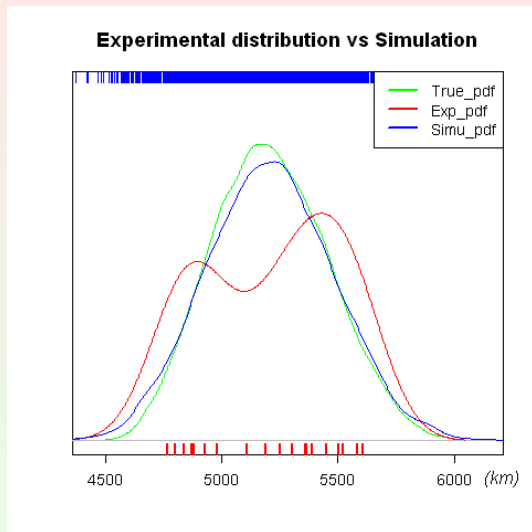
$$\hat{\theta}_{n, m}(h) = \underset{\theta \in \Theta}{\text{Argmin}} -\frac{1}{n} \sum_{i=1}^n \ln \left(\frac{1}{m} \sum_{j=1}^m K_{h_m}(Y_i^{\text{exp}} - Y_j^{\text{sim}}) \right)$$

- **A posteriori training data** : $(Y_1^{\text{exp}}, \dots, Y_n^{\text{exp}}, \hat{Y}_1^{\text{sim}}, \dots, \hat{Y}_m^{\text{sim}}) \rightarrow n + m = 3020!$

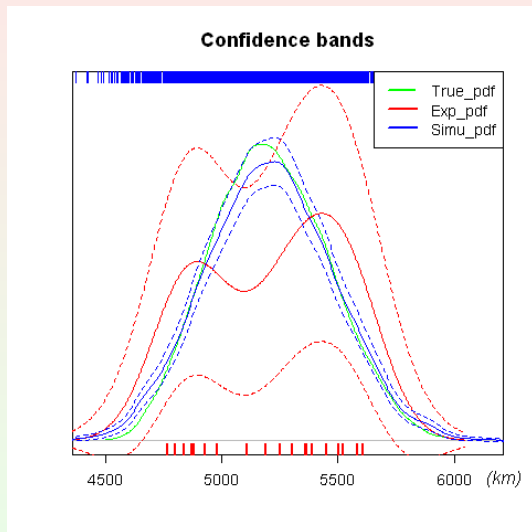
- QoI with *a priori* (Experimental) and *a posteriori* (Experimental + Simulated) training data



- QoI with *a priori* (Experimental) and *a posteriori* (Experimental + Simulated) training data



- QoI with *a priori* (Experimental) and *a posteriori* (Experimental + Simulated) training data



Other example: *Mean* study

Suppose that $Y \sim h(X, \theta_0)$ with $\theta_0 \in \Theta$

Recall that $Y^{sim} = h(X, \theta)$ with $\theta \in \Theta$

Let $Qol = \mathbb{E}(Y) = \varphi_Y = \varphi_{h, \theta_0}$

- $\varphi_{h, \theta} = \mathbb{E}_X(h(X, \theta))$, $\mathcal{D}(\varphi_{h, \theta}, \varphi_{h, \theta_0}) = \mathbb{E}(\varphi_{h, \theta} - \varphi_{h, \theta_0})^2$ (quadratic risk)

- let $\varphi(\mathcal{Y}_n^{exp}) = \frac{1}{n} \sum_{i=1}^n Y_i^{exp}$ and

$$\varphi(\mathcal{Y}_n^{exp}, \mathcal{Y}_m^{sim}) := \frac{1}{n+m} (Y_1^{exp} + \dots + Y_n^{exp} + \underbrace{Y_1^{sim} + \dots + Y_m^{sim}}_{\text{depend on } \theta !})$$

- **Question:** Is $\varphi(\mathcal{Y}_n^{exp}, \mathcal{Y}_m^{sim})$ "better" than $\varphi(\mathcal{Y}_n^{exp})$?

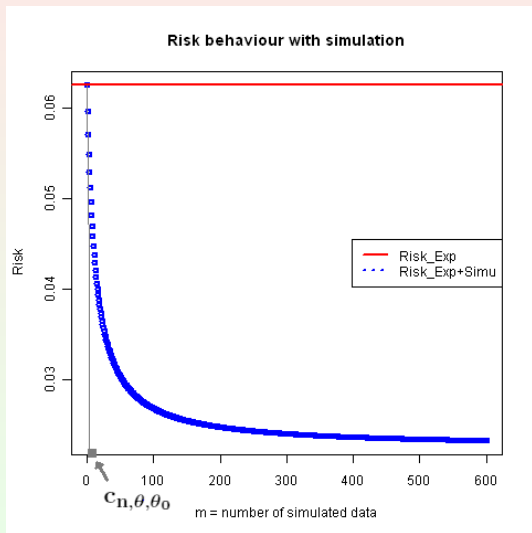
\Rightarrow turns out to have $\mathbf{R}_{n,m}(\theta) = \mathbb{E}(\varphi(\mathcal{Y}_n^{exp}, \mathcal{Y}_m^{sim}) - \varphi_{h, \theta_0})^2 \leq \mathbf{R}_n = \mathbb{E}(\varphi(\mathcal{Y}_n^{exp}) - \varphi_{h, \theta_0})^2$

Lemma

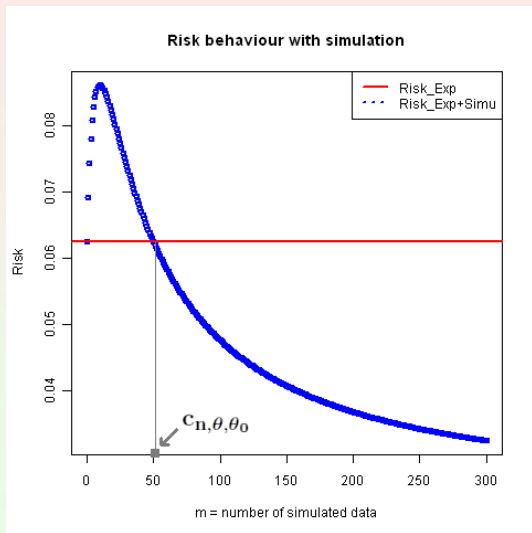
Let $n \in \mathbb{N}$. We show that $\exists \Theta_{sim}(n, \theta_0) \subset \Theta$ and $\exists c_{n, \theta, \theta_0} \in \mathbb{N}$ such that for all $\theta \in \Theta_{sim}(n, \theta_0)$

$$\mathbf{R}_{n,m}(\theta) \leq \mathbf{R}_n \quad \text{for all } m \geq c_{n, \theta, \theta_0}$$

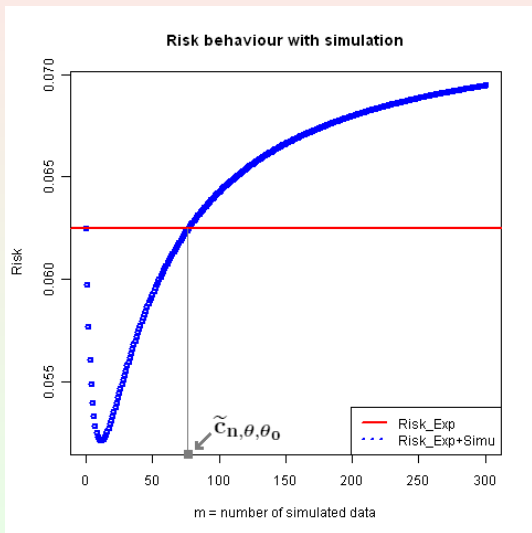
Mapping of $m \mapsto \mathbf{R}_{n,m}(\theta)$ and $m \mapsto \mathbf{R}_n$
For $\Rightarrow \theta$ in $\Theta_{sim}(n, \theta_0)$



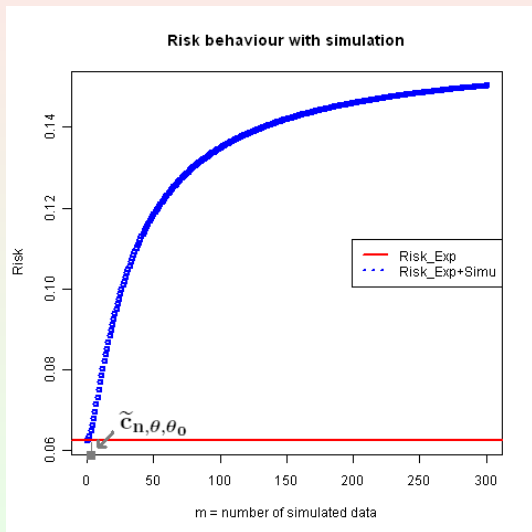
Mapping of $m \mapsto \mathbf{R}_{n,m}(\theta)$ and $m \mapsto \mathbf{R}_n$
For $\Rightarrow \theta$ in $\Theta_{sim}(n, \theta_0)$



Mapping of $m \mapsto \mathbf{R}_{n,m}(\theta)$ and $m \mapsto \mathbf{R}_n$
For $\Rightarrow \theta$ not in $\Theta_{sim}(n, \theta_0)$



Mapping of $m \mapsto \mathbf{R}_{n,m}(\theta)$ and $m \mapsto \mathbf{R}_n$
For $\Rightarrow \theta$ not in $\Theta_{sim}(n, \theta_0)$



Simulate or not to Simulate ?

- The only θ we dispose is $\hat{\theta}_{n,m}$
- **Question:** Does $\hat{\theta}_{n,m}$ belong to $\Theta_{sim}(n, \theta_0)$?
- Need of Central Limit Theorem !

- Rate of Convergence of the calibration procedure: fonction of n and m ...
 - Impact of Experimental and Simulated data on the estimation
 - For a given *Quantity of Interest* \Rightarrow how many n ? and how many m ?
 - etc ...

- Asymptotic Normality
 - Statistic studies
 - Confidence bands
 - etc...

- Sensitivity of this uncertainty analysis in relation to the *a priori* distribution of X .

- Robustness study: influence of the *QoI* on the *Model Selection*.

Thank you for your attention !