

Traitement des incertitudes

Josselin Garnier (Ecole Polytechnique)

- Problème général :

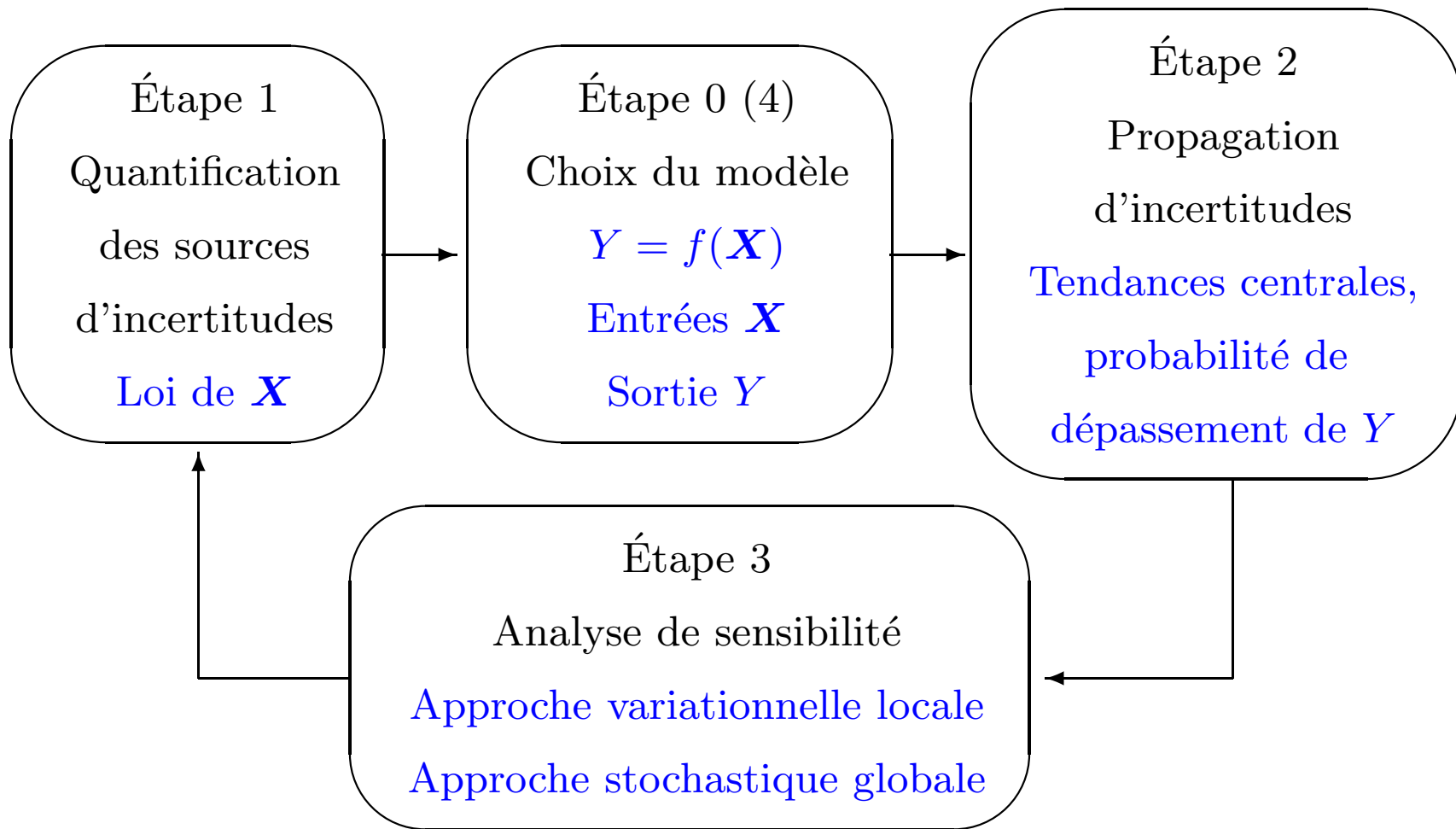
Comment modéliser les incertitudes et leur propagation dans les modèles physiques ou numériques ?

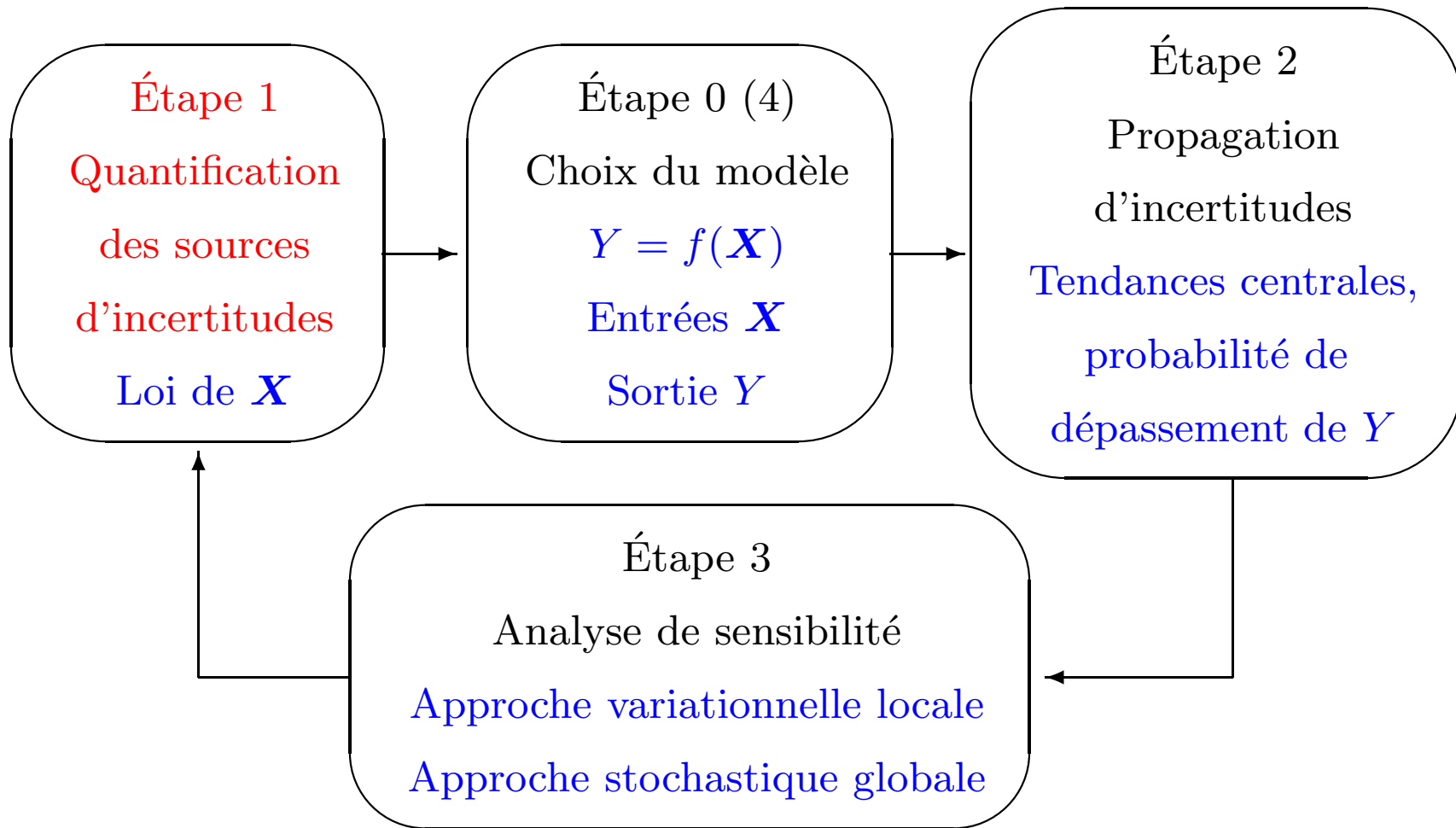
Comment estimer (quantifier) la dispersion de la sortie d'un code ou d'une expérience en fonction de la dispersion des paramètres d'entrée ?

Comment estimer (quantifier) la sensibilité de la sortie d'un code ou d'une expérience vis-à-vis d'un paramètre d'entrée particulier ?

Comment construire un métamodèle de la sortie d'un code ?

- But de l'exposé : faire un rapide état de l'art des différentes approches.





Partie I. Les sources d'incertitudes

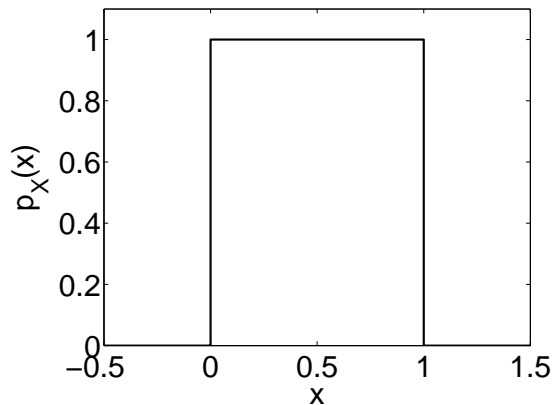
- Deux types d'entrée "incertaines" :
 - variables stochastiques : ces variables ont une variabilité naturelle résultant de phénomènes aléatoires (typiquement, une quantité soumise à des fluctuations dans un procédé de fabrication).
 - variables épistémiques : ces variables possèdent une valeur mais elle nous est inconnue, à cause d'un manque de connaissance (typiquement, une constante d'une loi physique).

- La modélisation par des lois de probabilités : les variables d'entrées sont traitées comme des variables aléatoires, de lois de probabilités données.

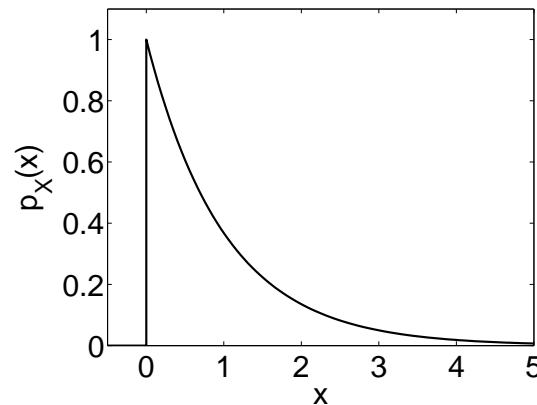
- La loi de probabilité d'une variable aléatoire réelle X caractérise la probabilité $\mathbb{P}(X \in [a, b])$ pour tout $a < b$.

Dans le cas d'une variable aléatoire continue (prend ses valeurs sur \mathbb{R} ou un intervalle de \mathbb{R}) : la loi est donnée par la densité $(p(x))_{x \in \mathbb{R}}$

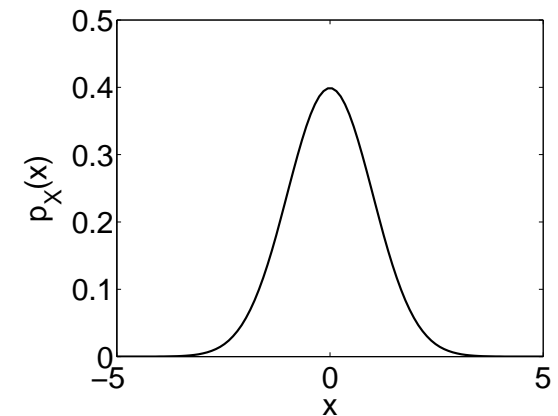
$$\mathbb{P}(X \in [a, b]) = \int_a^b p(x) dx, \quad \mathbb{P}(X \in [x, x + \delta x]) \simeq p(x) \delta x$$



uniforme



exponentielle



gaussienne

- Les variables d'entrées sont traitées comme des variables aléatoires, de lois de probabilités données.
- Pour déterminer la densité de la loi d'une variable aléatoire :
 - méthodes non-paramétriques à noyaux,
 - méthodes paramétriques d'ajustement à une loi analytique,
 - méthodes entropiques.
- Extensions : modèles hiérarchiques (classiques en analyse bayésienne).
- Remarque : Autres modélisations possibles :
 - distributions de possibilité (aussi appelées intervalles flous),
 - intervalles aléatoires utilisant les fonctions de croyance de Dempster-Shafer.

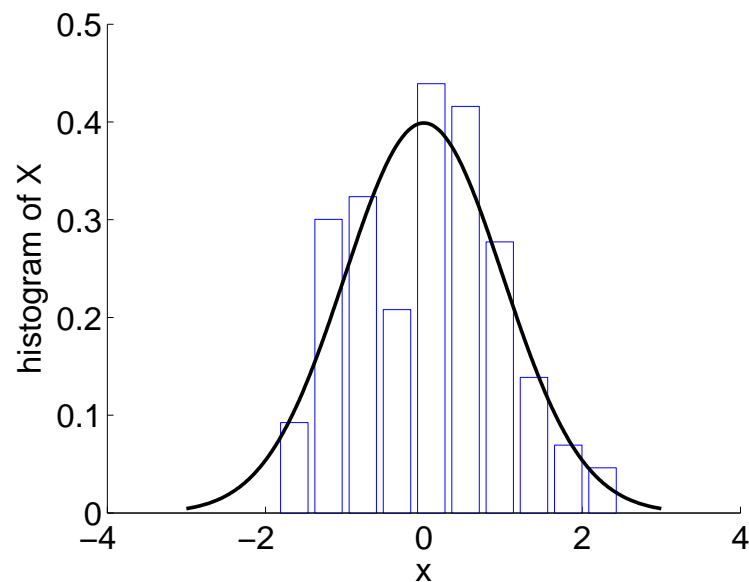
Cf: T. Hastie, R. Tibshirani et J. Friedman, *The Elements of Statistical Learning*, Springer, 2001.

- **Méthodes à noyaux.** Aussi appelée méthode de Parzen-Rozenblatt.
- Méthode non-paramétrique d'estimation de la densité d'une variable aléatoire.
 - se base sur un échantillon $(x_i)_{i=1,\dots,n}$ de réalisations indépendantes de la variable (on suppose la variable d'entrée de dimension un).
 - nécessite un nombre suffisant de données.
 - permet d'estimer la densité en tout point du support.
 - généralise la méthode d'estimation par un histogramme.
- Formule

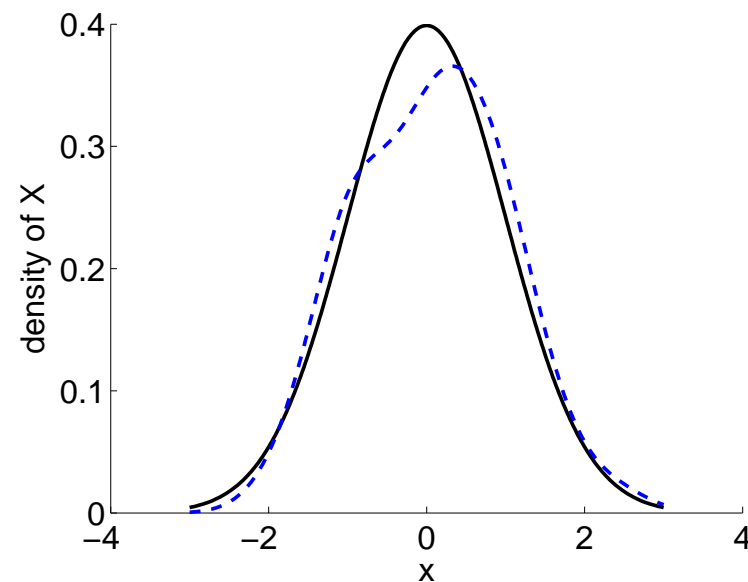
$$\hat{p}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

K est le noyau (souvent gaussien), h la fenêtre (régit le niveau de lissage).

- Exemple. Soit un échantillon de 200 données tirées selon une loi gaussienne :



histogramme



méthode à noyaux

- Résultats :
 - Il n'existe pas d'estimateur non-paramétrique qui converge plus vite que l'estimateur à noyaux (avec le choix de la fenêtre $h \simeq n^{-1/5}$).
 - La vitesse de convergence ($n^{-4/5}$ pour le risque quadratique) est plus faible que la vitesse typique des méthodes paramétriques (n^{-1} , voire plus rapide encore).

- Méthodes d'ajustement à une loi analytique.

Etant donné un échantillon $(x_i)_{i=1,\dots,n}$ de réalisations indépendantes de la variable.

On cherche à ajuster les paramètres d'une famille paramétrique de lois de probabilités aux données.

Exemple : famille de loi gaussiennes de densités de probabilités

$$p_{\mu,\sigma}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

(gaussienne de moyenne μ et de variance σ^2).

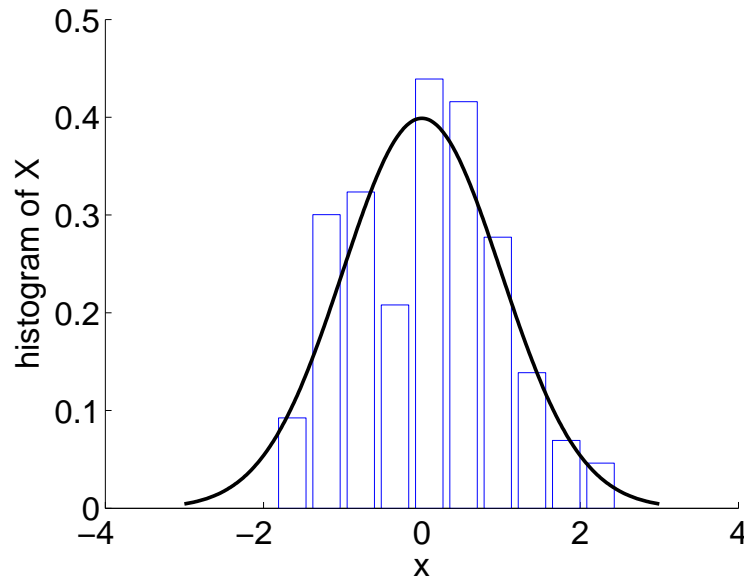
Nécessite un peu moins de données $(x_i)_{i=1,\dots,n}$ que la méthode à noyaux.

- Méthode des moments : on ajuste les paramètres de la loi analytique pour que ses premiers moments correspondent aux moments empiriques de l'échantillon.

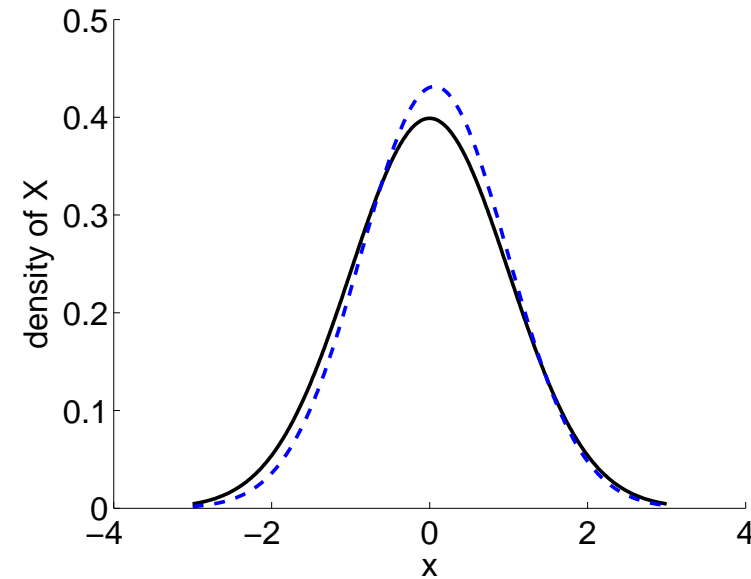
Exemple :

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Exemple. Soit un échantillon de 200 données tirées selon une loi gaussienne :



histogramme



méthode des moments

Résultat : risque quadratique converge en $1/n$ si la famille paramétrique contient la loi originale.

- Méthode du maximum de vraisemblance : on ajuste les paramètres $\boldsymbol{\theta}$ de la loi analytique pour que la log-vraisemblance soit maximale (en $\boldsymbol{\theta}$):

$$L_{\mathbf{x}}(\boldsymbol{\theta}) = \prod_{i=1}^n p_{\boldsymbol{\theta}}(x_i), \quad \ln L_{\mathbf{x}}(\boldsymbol{\theta}) = \sum_{i=1}^n \ln p_{\boldsymbol{\theta}}(x_i)$$

- Il y a un théorème de Bayes la-dessous : $L_{\mathbf{x}}(\boldsymbol{\theta})$ est la vraisemblance des données $\mathbf{x} = (x_i)_{i=1, \dots, n}$ sachant le paramètre $\boldsymbol{\theta}$, et est aussi la vraisemblance du paramètre $\boldsymbol{\theta}$ sachant les données \mathbf{x} sans a priori sur $\boldsymbol{\theta}$.
- L'estimateur du maximum de vraisemblance peut exister et être unique, ne pas être unique, ou ne pas exister.
- Dans les bons cas, il converge, son risque quadratique est en $1/n$, il est asymptotiquement normal, asymptotiquement efficace (i.e., on ne peut pas faire mieux, asymptotiquement).
- Exemple : $\boldsymbol{\theta} = (\mu, \sigma^2)$:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

- Possibilité de prendre en compte un a priori (avis d'expert ou étude antérieure) sur $\boldsymbol{\theta}$ (\rightarrow méthode du maximum a posteriori).

- Méthode bayésienne : fournit une distribution sur les paramètres de la loi (et pas seulement une estimation); permet d'incorporer une idée a priori sur les paramètres; on obtient un modèle hiérarchique.

Exemple : $\mathbf{x} = (x_k)_{k=1,\dots,n}$ échantillon de loi $\mathcal{N}(\mu, \sigma_{\text{mes}}^2)$ avec μ inconnu et σ_{mes} connu :

$$p(\mathbf{x}|\mu) = \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma_{\text{mes}}} \exp\left(-\frac{(x_k - \mu)^2}{2\sigma_{\text{mes}}^2}\right)$$

$$\hookrightarrow p_{\text{post}}(\mu|\mathbf{x}) \approx \exp\left(-\sum_{k=1}^n \frac{(x_k - \mu)^2}{2\sigma_{\text{mes}}^2}\right) \approx \exp\left(-\frac{(\mu - \frac{1}{n} \sum_{k=1}^n x_k)^2}{2\frac{\sigma_{\text{mes}}^2}{n}}\right)$$

la loi a posteriori de μ , étant données les observations \mathbf{x} , est

$$\mathcal{N}\left(\frac{1}{n} \sum_{k=1}^n x_k, \frac{\sigma_{\text{mes}}^2}{n}\right)$$

A posteriori :

$$\mathbb{E}[f(X)] = \int \mathbb{E}_{\mathcal{N}(\mu, \sigma_{\text{mes}}^2)}[f(X)] \frac{1}{\sqrt{2\pi \frac{\sigma_{\text{mes}}^2}{n}}} \exp\left(-\frac{(\mu - \frac{1}{n} \sum_{k=1}^n x_k)^2}{2\frac{\sigma_{\text{mes}}^2}{n}}\right) d\mu$$

- Test de qualité d'ajustement permet de dire si les données $(x_i)_{i=1,\dots,n}$ sont compatibles avec la famille paramétrique p_{θ} .

Attention avec les tests : la réponse du test est "je ne peux pas rejeter l'hypothèse que les données sont compatibles avec la famille de lois" ou "je rejette l'hypothèse que les données sont compatibles avec la famille de lois".

- Le choix de la famille est un jugement d'expert; peut-être motivée par des considérations théoriques, un théorème central limite pour justifier une gaussienne, un principe d'entropie, ...

- Méthode du maximum d'entropie.

On choisit la loi (densité $p(x)$) qui maximise l'entropie (le manque d'information)

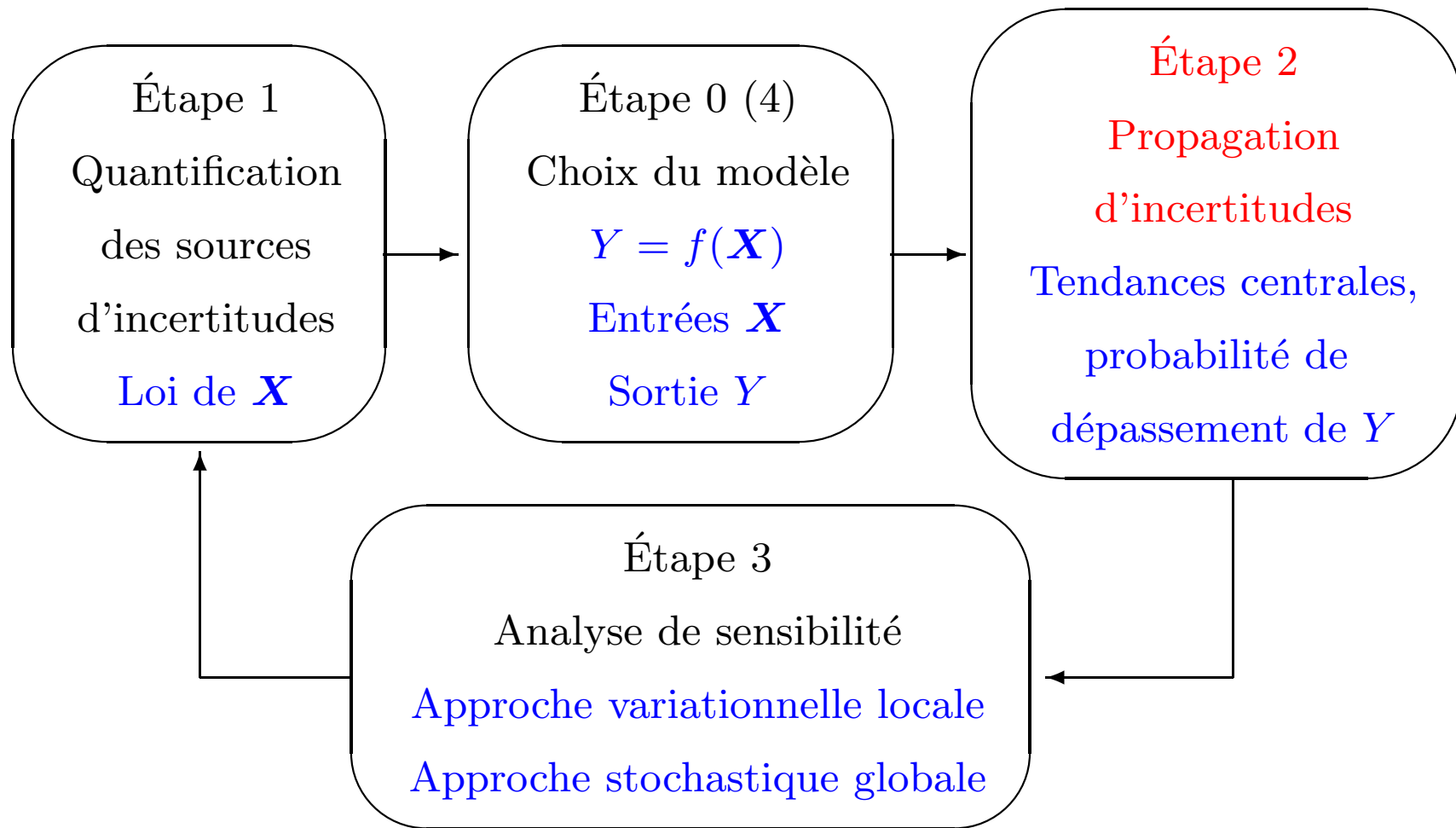
$$E(p) = - \int p(x) \ln p(x) dx$$

étant données certaines informations (moyenne, variance, intervalles de définition).

Exemples :

- la loi qui maximise l'entropie d'une variable aléatoire de moyenne et de variance prescrites est une loi gaussienne.
- la loi qui maximise l'entropie d'une variable aléatoire à valeurs dans un intervalle prescrit est une loi uniforme.
- la loi qui maximise l'entropie d'une variable aléatoire à valeurs positives et à moyenne prescrite est une loi exponentielle.

Joue un rôle dans les méthodes bayésiennes pour définir les lois a priori.



Partie II. Propagation d'incertitudes

- Contexte : code de calcul informatique ou expérience modélisé par

$$Y = f(\mathbf{X})$$

avec Y =variable de sortie

$\mathbf{X} = (X_i)_{i=1,\dots,d}$ variables d'entrée

f = boîte noire déterministe

On se donne : la loi de \mathbf{X} .

- But : estimation d'une quantité

$$\mathbb{E}[\psi(Y)]$$

avec une barre d'erreur et le minimum de simulations/expériences.

- Exemples (pour Y à valeurs réelles) :

$\psi(y) = y \rightarrow$ moyenne de Y , i.e. $\mathbb{E}[Y]$

$\psi(y) = y^2 \rightarrow$ variance de Y , i.e. $\text{Var}(Y) = \mathbb{E}[(Y - \mathbb{E}[Y])^2] = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$

$\psi(y) = \mathbf{1}_{[a,\infty)}(y) \rightarrow$ probabilité de dépasser le seuil a , i.e. $\mathbb{P}(Y \geq a)$.

Méthode analytique

- La quantité à estimer est une intégrale d -dimensionnelle :

$$I = \mathbb{E}[\psi(Y)] = \mathbb{E}[F(\mathbf{X})] = \int_{\mathbb{R}^d} F(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

où $p(\mathbf{x})$ est la densité de probabilité de \mathbf{X} et $F(\mathbf{x}) = \psi(f(\mathbf{x}))$.

Dans des cas simples (quand la densité p et la fonction F sont connues de manière analytique), on peut parfois évaluer l'intégrale de manière exacte.

Situation "exceptionnelle".

Méthodes de quadrature

- La quantité à estimer est une intégrale d -dimensionnelle :

$$I = \mathbb{E}[\psi(Y)] = \mathbb{E}[F(\mathbf{X})] = \int_{\mathbb{R}^d} F(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

où $p(\mathbf{x})$ est la densité de probabilité de \mathbf{X} et $F(\mathbf{x}) = \psi(f(\mathbf{x}))$.

Si $p(\mathbf{x}) = \prod_{i=1}^d p_0(x_i)$, alors on peut appliquer une méthode de quadrature gaussienne avec grille pleine de n^d points :

$$\hat{I} = \sum_{j_1=1}^n \cdots \sum_{j_d=1}^n \rho_{j_1} \cdots \rho_{j_d} F(\xi_{j_1}, \dots, \xi_{j_d})$$

avec $(\rho_j)_{j=1, \dots, n}$ poids et $(\xi_j)_{j=1, \dots, n}$ nœuds de quadrature gaussienne associée à la densité p_0 .

Il existe aussi des méthodes de quadrature avec des grilles creuses (Smolyak).

- Les méthodes de quadrature sont efficaces avec :
 - des conditions de régularité sur $\mathbf{x} \rightarrow F(\mathbf{x})$ (et pas seulement f),
 - une dimension d petite (même avec des grilles creuses de type Smolyak).

Elles réclament sinon beaucoup d'appels au code.

Méthode de Monte Carlo

On tire un n -échantillon $(\mathbf{X}^{(k)})_{k=1,\dots,n}$ (des réalisations indépendantes de \mathbf{X}).

Estimateur de $I = \mathbb{E}[F(\mathbf{X})]$:

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)})$$

Estimation non-biaisée :

$$\mathbb{E}[\hat{I}_n] = I \quad \text{pour tout } n$$

Convergence :

$$\hat{I}_n \xrightarrow[n \rightarrow \infty]{} I \quad \text{avec probabilité 1}$$

Erreur (risque quadratique) :

$$\mathbb{E}[(\hat{I}_n - I)^2] = \text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}(F(\mathbf{X}))$$

Erreur relative :

$$\frac{\mathbb{E}[(\hat{I}_n - I)^2]^{1/2}}{I} = \frac{1}{\sqrt{n}} \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\mathbb{E}[F(\mathbf{X})]}$$

Caractéristiques :

- 1) possibilité d'obtenir des intervalles de confiance,
- 2) pas de régularité requise sur F (avec $F(\mathbf{x}) = \psi(f(\mathbf{x}))$),
- 3) vitesse de convergence indépendante de la dimension (mais lente).

Intervalles de confiance

Question : A partir de l'échantillon $(\mathbf{X}^{(k)})_{k=1,\dots,n}$, l'estimateur \hat{I}_n donne une valeur approchée de I , d'autant meilleure que n est grand. Comment quantifier précisément l'erreur ?

Réponse : On construit un intervalle de confiance au niveau 0.95, i.e. un intervalle $[\hat{a}_n, \hat{b}_n]$ tel que

$$\mathbb{P} \left(I \in [\hat{a}_n, \hat{b}_n] \right) \geq 0.95$$

Construction basée sur le **théorème central limite** :

$$\mathbb{P} \left(\left| \hat{I}_n - I \right| < c \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\sqrt{n}} \right) \xrightarrow{n \rightarrow \infty} \frac{2}{\sqrt{2\pi}} \int_0^c e^{-x^2/2} dx$$

Le membre de droite vaut 0.95 si $c = 1.96$. Donc

$$\mathbb{P} \left(I \in \left[\hat{I}_n - 1.96 \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\sqrt{n}}, \hat{I}_n + 1.96 \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\sqrt{n}} \right] \right) \simeq 0.95$$

$$\mathbb{P} \left(I \in \left[\hat{I}_n - 1.96 \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\sqrt{n}}, \hat{I}_n + 1.96 \frac{\text{Var}(F(\mathbf{X}))^{1/2}}{\sqrt{n}} \right] \right) \simeq 0.95$$

Les bornes sont encore inconnues car on ne connaît pas $\text{Var}(F(\mathbf{X}))$! Deux solutions :

- solution conservatrice, du type $\text{Var}(F(\mathbf{X})) \leq \|F\|_\infty^2$, et alors

$$\mathbb{P} \left(I \in \left[\hat{I}_n - 1.96 \frac{\|F\|_\infty}{\sqrt{n}}, \hat{I}_n + 1.96 \frac{\|F\|_\infty}{\sqrt{n}} \right] \right) \geq 0.95$$

- asymptotiquement, on remplace $\text{Var}(F(\mathbf{X}))$ dans les bornes par son estimateur empirique $\hat{\sigma}_n^2$:

$$\mathbb{P} \left(I \in \left[\hat{I}_n - 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}}, \hat{I}_n + 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}} \right] \right) \simeq 0.95$$

où

$$\hat{\sigma}_n = \left(\frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)})^2 - \hat{I}_n^2 \right)^{1/2}$$

Conclusion : Il n'y a aucun intervalle borné de \mathbb{R} dont on puisse dire avec certitude qu'il contient I , mais il y a des intervalles, dits intervalle de confiance, dont on peut dire qu'ils contiennent I avec une probabilité proche de 1.

Propagation d'incertitudes par métamodèles

On remplace f par un métamodèle (modèle réduit) f_r et on applique une des techniques précédentes (analytique, quadrature, Monte Carlo).

→ On peut faire beaucoup d'appels au métamodèle.

→ Le choix du métamodèle est critique.

→ Le contrôle de l'erreur n'est pas simple.

Développements de Taylor

- On approche la sortie $Y = f(\mathbf{X})$ par un développement de Taylor $Y_r = f_r(\mathbf{X})$.
- Exemple :
 - On souhaite estimer $\mathbb{E}[Y]$ et $\text{Var}(Y)$ pour $Y = f(\mathbf{X})$ avec X_i décorréliées, $\mathbb{E}[X_i] = \mu_i$ et $\text{Var}(X_i) = \sigma_i^2$ connus, σ_i^2 petits.
 - On approche $Y = f(\mathbf{X})$ par $Y_r = f_r(\mathbf{X}) = f(\boldsymbol{\mu}) + \nabla f(\boldsymbol{\mu}) \cdot (\mathbf{X} - \boldsymbol{\mu})$. On trouve :

$$\mathbb{E}[Y] \simeq \mathbb{E}[Y_r] = f(\boldsymbol{\mu}), \quad \text{Var}(Y) \simeq \text{Var}(Y_r) = \sum_{i=1}^d \partial_{x_i} f(\boldsymbol{\mu})^2 \sigma_i^2$$

On a juste besoin de calculer $f(\boldsymbol{\mu})$ et $\nabla f(\boldsymbol{\mu})$ (calcul du gradient par différences finies ou par différenciation automatique, donc en gros $d + 1$ appels à f).

- Rapide, analytique, permet de calculer approximativement des tendances centrales de la sortie (moyenne, variance).
- Convenable pour des petites variations des paramètres d'entrée et un modèle régulier (qu'on peut linéariser).
- Approche "locale". En général, pas de contrôle de l'erreur.

FORM-SORM

$$P = \mathbb{P}(f(\mathbf{X}) \geq a) = \mathbb{P}(\mathbf{X} \in \mathcal{F}) = \int_{\mathcal{F}} p(\mathbf{x})d\mathbf{x}, \quad \mathcal{F} = \{\mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \geq a\}$$

- Méthode FORM-SORM, analytique mais approchée, sans contrôle d'erreur.
- on suppose que les X_i sont indépendants et de loi gaussienne de moyenne zéro et de variance un (ou on se ramène à ce cas par transformation isoprobabiliste).
- on trouve par optimisation (sous contrainte) le point \mathbf{x}_a de dépassement de seuil (i.e. $f(\mathbf{x}_a) = a$) le plus proche de l'origine.
- on approche la surface de défaillance $\{\mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) = a\}$ par une surface régulière $\hat{\mathcal{F}}$ qui permette de faire un calcul analytique $\hat{P} = \int_{\hat{\mathcal{F}}} p(\mathbf{x})d\mathbf{x}$:

- un hyperplan pour FORM

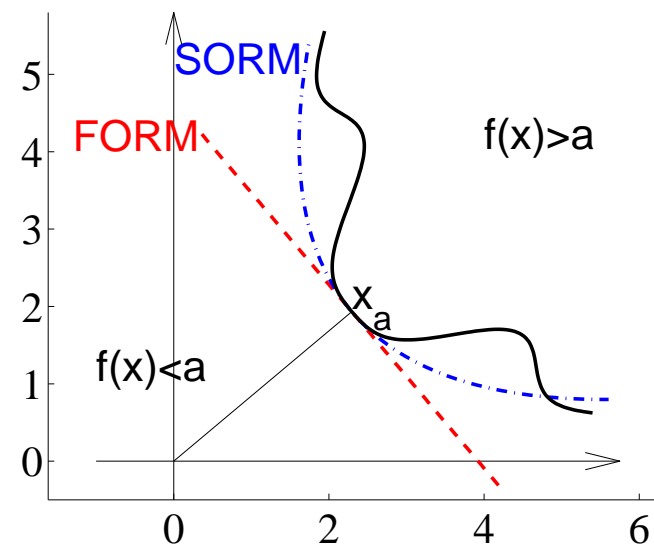
(et alors $\hat{P} = \frac{1}{2}\text{erfc}\left(\frac{|\mathbf{x}_a|}{\sqrt{2}}\right)$),

- une forme quadratique pour SORM

(et alors $\hat{P} =$ formule de Breitung).

Cf: O. Ditlevsen et H.O. Madsen,

Structural reliability methods, Wiley, 1996.



Techniques de réduction de variance

Objectif : réduire la variance de l'estimateur de Monte Carlo classique :

$$\mathbb{E}[(\hat{I}_n - I)^2] = \frac{1}{n} \text{Var}(F(\mathbf{X}))$$

Les méthodes

- Echantillonnage préférentiel
- Variables de contrôle
- Variables antithétiques
- Stratification

ont pour but de réduire la constante, en restant proches de l'esprit Monte Carlo (parallélisable).

Les méthodes

- Quasi-Monte Carlo

ont pour but de changer le $1/n$.

Les méthodes

- Systèmes de particules en interaction

s'éloignent de l'esprit Monte Carlo (séquentielle).

Echantillonnage préférentiel (importance sampling)

- Observation : la représentation de I comme espérance n'est pas unique. Si \mathbf{X} est à densité p :

$$I = \mathbb{E}_p[F(\mathbf{X})] = \int F(\mathbf{x})p(\mathbf{x})d\mathbf{x} = \int \frac{F(\mathbf{x})p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} = \mathbb{E}_q\left[\frac{F(\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})}\right]$$

L'utilisateur est libre dans le choix de la densité q .

- Idée : Dans la situation où on sait que $F(\mathbf{X})$ est surtout sensible à certaines valeurs de \mathbf{X} , au lieu de tirer les $\mathbf{X}^{(k)}$ selon la densité originale $p(\mathbf{x})$ de \mathbf{X} , on les tire selon une densité "biaisée" $q(\mathbf{x})$ qui favorise les valeurs de \mathbf{X} dans la zone d'importance.

- En considérant la représentation

$$I = \mathbb{E}_p[F(\mathbf{X})] = \mathbb{E}_q\left[F(\mathbf{X})\frac{p(\mathbf{X})}{q(\mathbf{X})}\right]$$

on propose l'estimateur :

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)}) \frac{p(\mathbf{X}^{(k)})}{q(\mathbf{X}^{(k)})}$$

où $(\mathbf{X}^{(k)})_{k=1,\dots,n}$ est un n -échantillon tiré selon q .

- Estimateur non-biaisé : $\mathbb{E}_q[\hat{I}_n] = I$.

- Estimateur convergent :

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)}) \frac{p(\mathbf{X}^{(k)})}{q(\mathbf{X}^{(k)})} \xrightarrow{n \rightarrow \infty} \mathbb{E}_q \left[F(\mathbf{X}) \frac{p(\mathbf{X})}{q(\mathbf{X})} \right] = \mathbb{E}_p [F(\mathbf{X})] = I$$

- Variance de l'estimateur :

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}_q \left(F(\mathbf{X}) \frac{p(\mathbf{X})}{q(\mathbf{X})} \right) = \frac{1}{n} \left(\mathbb{E}_p \left[F(\mathbf{X})^2 \frac{p(\mathbf{X})}{q(\mathbf{X})} \right] - \mathbb{E}_p [F(\mathbf{X})]^2 \right)$$

- En choisissant bien q , on peut fortement réduire la variance. En fait, en choisissant (en supposant $F \geq 0$)

$$q_{opt}(\mathbf{x}) = \frac{F(\mathbf{x})p(\mathbf{x})}{I}$$

on trouve

$$\text{Var}(\hat{I}_n) = 0$$

Mais $q_{opt}(\mathbf{x})$ dépend de I ! (\hookrightarrow méthodes séquentielles/adaptatives).

- Points pratiques importants pour pouvoir implémenter la méthode :
 - Il faut savoir simuler \mathbf{X} de loi de densité q .
 - Il faut savoir calculer le rapport de vraisemblance $\frac{p(\mathbf{x})}{q(\mathbf{x})}$.

- Exemple : On veut estimer

$$I = \mathbb{E}[F(X)]$$

avec $X \sim \mathcal{N}(0, 1)$ et $F(x) = \mathbf{1}_{[4, \infty[}(x)$.

$$I = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \mathbf{1}_{[4, \infty[}(x) e^{-\frac{x^2}{2}} dx = \frac{1}{2} \operatorname{erfc}\left(\frac{4}{\sqrt{2}}\right) \simeq 3.17 \cdot 10^{-5}$$

Monte Carlo :

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_k \geq 4}, \quad X_k \sim \mathcal{N}(0, 1)$$

On a $\operatorname{Var}(\hat{I}_n) = \frac{1}{n} 3.17 \cdot 10^{-5}$.

Echantillonnage préférentiel : on tire X_k selon la loi $\mathcal{N}(4, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_k \geq 4} \frac{e^{-\frac{X_k^2}{2}}}{e^{-\frac{(X_k-4)^2}{2}}} = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_k \geq 4} e^{-4X_k+8}, \quad X_k \sim \mathcal{N}(4, 1)$$

On a $\operatorname{Var}(\hat{I}_n) = \frac{1}{n} 5.53 \cdot 10^{-8}$.

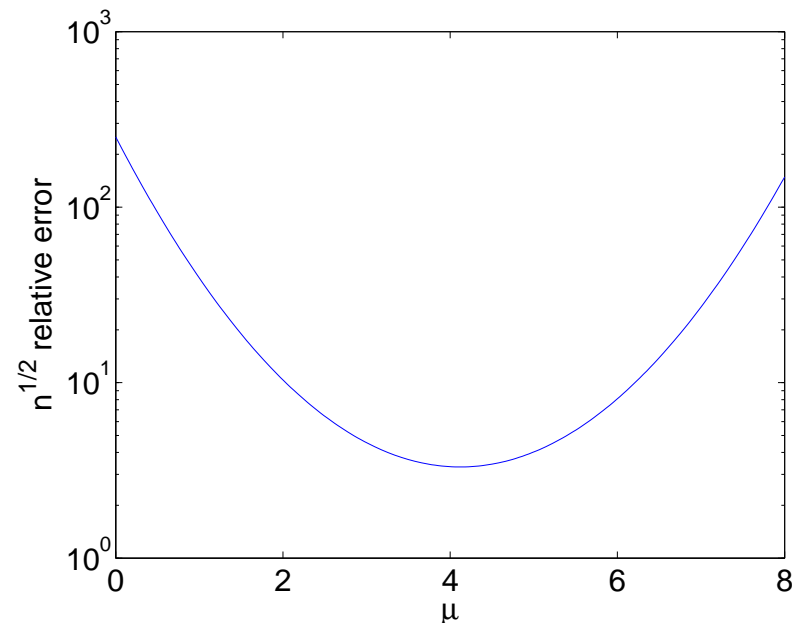
Il faut 1000 fois moins de simulations avec la méthode IS pour atteindre la même précision !

Attention cependant, il ne faut pas trop pousser.

Echantillonnage préférentiel : on tire X_k selon la loi $\mathcal{N}(\mu, 1)$.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_k \geq 4} \frac{e^{-\frac{X_k^2}{2}}}{e^{-\frac{(X_k - \mu)^2}{2}}} = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{X_k \geq 4} e^{-\mu X_k + \frac{\mu^2}{2}}, \quad X_k \sim \mathcal{N}(\mu, 1)$$

On a $\text{Var}(\hat{I}_n) = \frac{1}{n} \frac{e^{\mu^2}}{2} \text{erfc}\left(\frac{4+\mu}{\sqrt{2}}\right) - \frac{1}{n} I^2$, ce qui donne pour l'erreur relative normalisée $\sqrt{n} \mathbb{E}[(\hat{I}_n - I)^2]^{1/2} / I$:



Si on pousse trop, les fluctuations des rapports de vraisemblance deviennent trop grandes.

Variables de contrôle

- Rappel : Le but est d'estimer $I = \mathbb{E}[F(\mathbf{X})]$ pour \mathbf{X} un vecteur aléatoire et $F(\mathbf{x}) = \psi(f(\mathbf{x}))$ une fonction déterministe.
- Dans la situation où on dispose d'un modèle réduit $f_r(\mathbf{x})$.
- Méthode d'échantillonnage préférentiel : on calcule (on approche) la densité optimale $q_{opt,r}(\mathbf{x}) = \frac{\psi(f_r(\mathbf{x}))p(\mathbf{x})}{I_r}$, avec $I_r = \int \psi(f_r(\mathbf{x}))p(\mathbf{x})d\mathbf{x}$, puis on l'utilise comme densité d'importance.

- Méthode de variables de contrôle :

On note $F(\mathbf{x}) = \psi(f(\mathbf{x}))$, $F_r(\mathbf{x}) = \psi(f_r(\mathbf{x}))$.

On suppose qu'on connaît $I_r = \mathbb{E}[F_r(\mathbf{X})]$.

En considérant la représentation

$$I = \mathbb{E}[F(\mathbf{X})] = I_r + \mathbb{E}[F(\mathbf{X}) - F_r(\mathbf{X})]$$

on propose l'estimateur :

$$\hat{I}_n = I_r + \frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)}) - F_r(\mathbf{X}^{(k)}),$$

où $(\mathbf{X}^{(k)})_{k=1,\dots,n}$ est un n -échantillon (tiré selon p).

Estimateur :

$$\hat{I}_n = I_r + \frac{1}{n} \sum_{k=1}^n F(\mathbf{X}^{(k)}) - F_r(\mathbf{X}^{(k)})$$

Cet estimateur est sans biais et convergent.

Sa variance est :

$$\text{Var}(\hat{I}_n) = \frac{1}{n} \text{Var}(F(\mathbf{X}) - F_r(\mathbf{X}))$$

↪ Cette méthode peut réduire la variance.

- Exemple : on souhaite estimer

$$I = \mathbb{E}[f(X)]$$

avec $X \sim \mathcal{U}(0, 1)$, $f(x) = \exp(x)$.

Résultat : $I = e - 1 \simeq 1.72$.

Monte Carlo.

$$\hat{I}_n = \frac{1}{n} \sum_{k=1}^n \exp[X^{(k)}]$$

Variance de l'estimateur MC = $\frac{1}{n}(2e - 1) \simeq \frac{1}{n}4.44$.

Variable de contrôle. Modèle réduit : $f_r(x) = 1 + x$ (ici $I_r = \frac{3}{2}$). Estimateur VC :

$$\hat{I}_n = I_r + \frac{1}{n} \sum_{k=1}^n \left\{ \exp[X^{(k)}] - 1 - X^{(k)} \right\}$$

Variance de l'estimateur VC = $\frac{1}{n}(3e - \frac{e^2}{2} - \frac{53}{12}) \simeq \frac{1}{n}0.044$.

Il faut donc 100 fois moins de simulations avec l'estimateur VC !

- Application : Méthodes de Romberg statistiques pour l'estimation de

$$I = \mathbb{E}[\psi(f(\mathbf{X}))]$$

On dispose d'un code léger f_r en plus du code lourd f . Le rapport du coût calcul entre un appel à f et un appel à f_r est $q > 1$.

Estimateur

$$\hat{I}_n = \frac{1}{n_r} \sum_{i=1}^{n_r} F_r(\tilde{\mathbf{X}}^{(k)}) + \frac{1}{n} \sum_{i=1}^n F(\mathbf{X}^{(k)}) - F_r(\mathbf{X}^{(k)})$$

avec $n_r \gg n$, $F(\mathbf{x}) = \psi(f(\mathbf{x}))$, $F_r(\mathbf{x}) = \psi(f_r(\mathbf{x}))$.

Allocation entre appels au code lourd et appels au code léger à optimiser sous la contrainte $n_r/q + n(1 + 1/q) = n_{\text{tot}}$.

Compromis classique entre erreur d'approximation et erreur d'estimation.

- Utilisé dans le cas où $f(\mathbf{X})$ est la solution d'une équation différentielle discrétisée finement, avec $f_r(\mathbf{X})$ la solution avec un schéma de discrétisation grossier (MultiLevel Monte Carlo).

↔ Il ne faut pas utiliser le code avec la précision maximale pour tous les appels dans un problème de propagation d'incertitudes (ou d'analyse de sensibilité) !

Suite à discrédance faible (quasi Monte Carlo)

on tire l'échantillon de manière moins aléatoire que MC, pour combler les trous qui se forment naturellement dans un échantillon aléatoire.

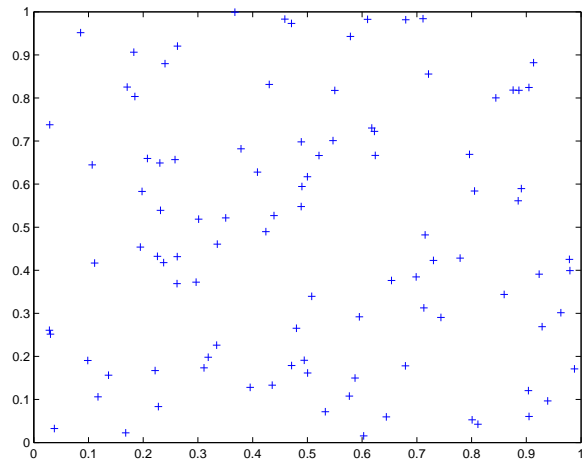
Cette technique

- réduit la variance si f a un peu de régularité et/ou de monotonie; on peut aller jusqu'à une variance en $C_d(\log n)^{s(d)}/n^2$,
- marche en dimension pas trop grande,
- représente un intermédiaire entre MC et quadrature usuelle,
- en compétition avec des méthodes de quadrature creuse (Smolyak).

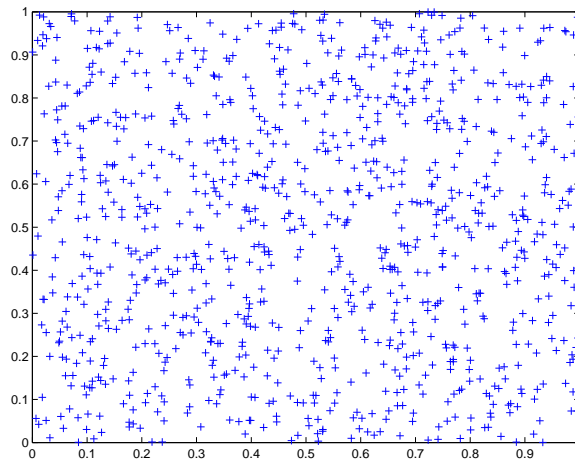
Attention :

- il n'est pas facile de rajouter des points,
- on n'a pas d'estimée d'erreur (sans hypothèse supplémentaire),
- la méthode n'est pas adaptée pour l'estimation de probabilité d'événement rare.

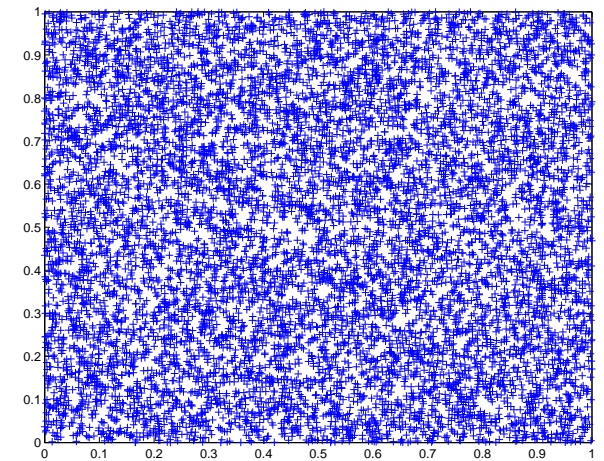
Exemple : échantillon Monte Carlo.



$n = 100$

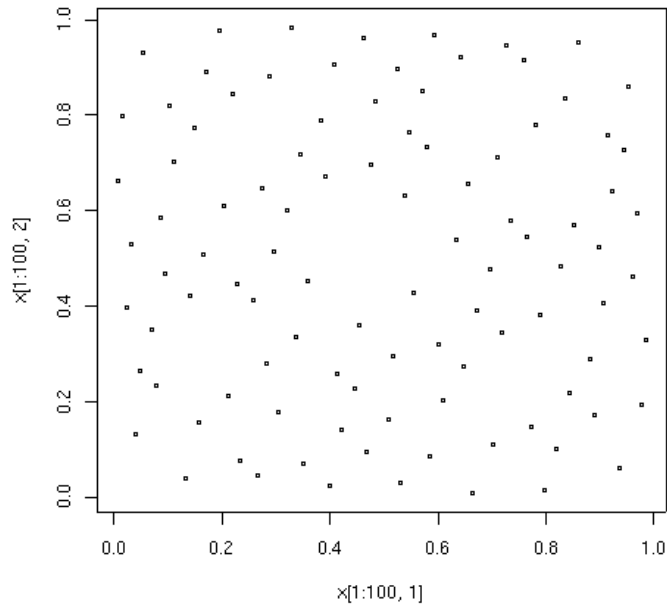


$n = 1000$

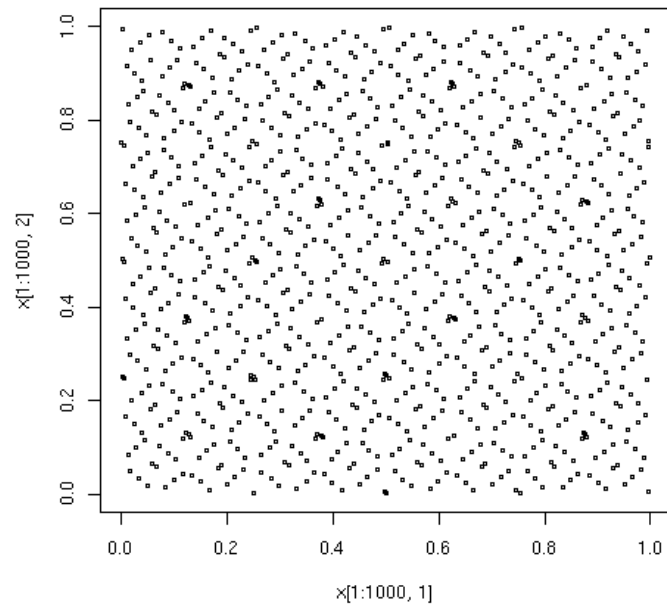


$n = 10000$

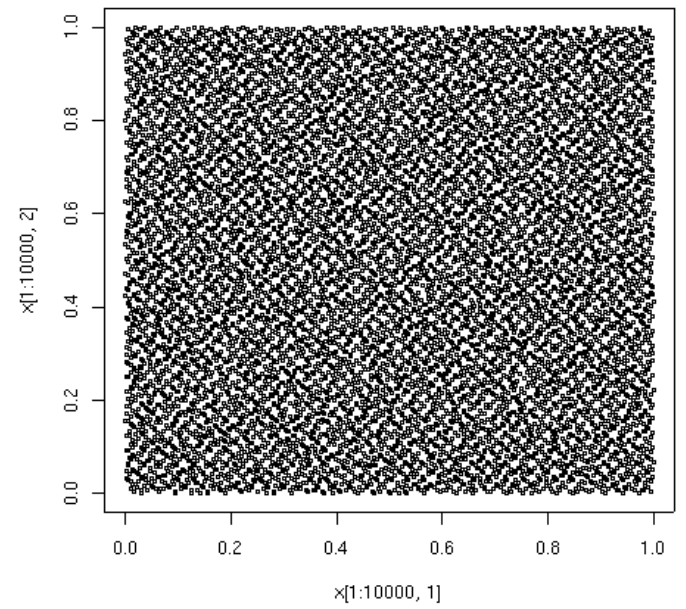
Exemple : suites de Sobol en dimension 2.



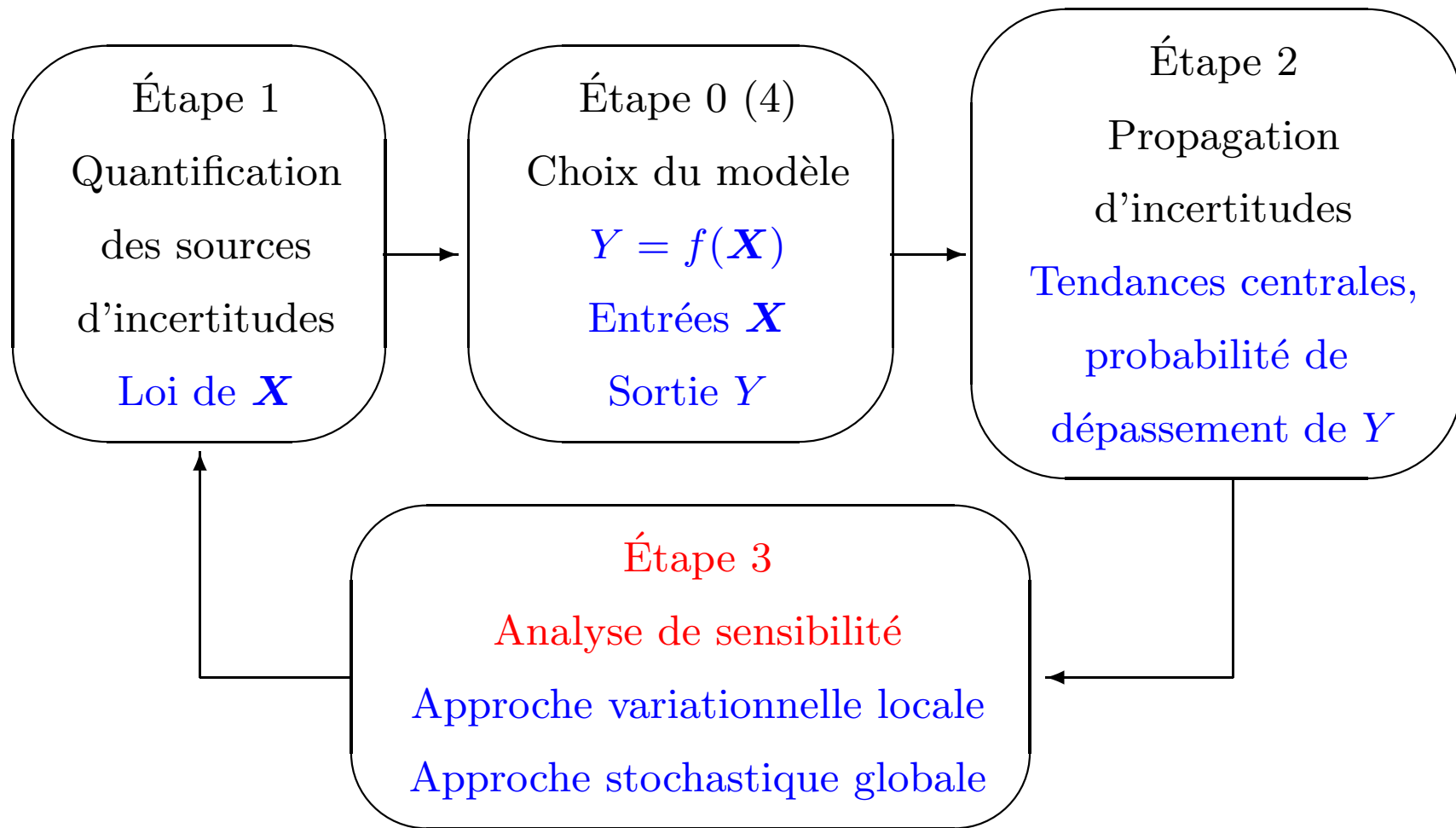
$n = 100$



$n = 1000$



$n = 10000$



Partie III. Analyse de sensibilité

- Contexte : code de calcul informatique ou expérience modélisé par

$$Y = f(\mathbf{X})$$

avec Y =variable de sortie réelle

$\mathbf{X} = (X_1, \dots, X_d)$ variables d'entrée

f = boîte noire déterministe

- But : expliquer la variabilité de la réponse Y en fonction des X_i .
- Objectifs principaux :
 - améliorer la compréhension du phénomène.
 - réduire l'incertitude d'un modèle, en identifiant les variables les plus influentes dans un domaine de valeurs de la sortie (\rightarrow il devient prioritaire de réduire la variabilité de ces entrées).
 - simplifier ou alléger le modèle, en fixant les variables les moins influentes.

On se donne : la loi des X_i (gaussienne, uniforme,...).

But de l'analyse de sensibilité :

- qualitative : identifier les paramètres importants parmi les nombreux paramètres (screening ou criblage).
- quantitative : déterminer la part de la **variance** de la sortie Y due à une variable d'entrée ou un sous-ensemble de variables d'entrée X_i .

Références :

A. Saltelli, K. Chan et E. M. Scott, Sensitivity analysis, Wiley, 2001

A. Saltelli et S. Tarantola, Sensivity analysis in practice, Wiley, 2004

Criblage

Un modèle comportant beaucoup de variables d'entrée est difficile à explorer.

Souvent, seulement quelques entrées sont influentes.

Objectif : identifier rapidement les quelques h entrées influentes parmi les d entrées, avec n calculs.

Méthodes possibles :

- méthode OAT (One factor At a Time) : calcul de gradient avec $n = d + 1$.
- méthode de Morris avec $n = R(d + 1)$.
- criblage par groupes, bifurcation séquentielle (avec des hypothèses de monotonie) avec $n < d$ et $h \ll d$.

- Méthode OAT.

- On fait $d + 1$ appels au code en \mathbf{x}_0 et $\mathbf{x}_0 + \delta \mathbf{e}_j$, avec \mathbf{x}_0 un point de référence, δ un pas, et $(\mathbf{e}_j)_{j=1,\dots,d}$ base canonique de \mathbb{R}^d .

- Indice de sensibilité :

$$S_j^{\mathbf{x}_0} = \frac{f(\mathbf{x}_0 + \delta \mathbf{e}_j) - f(\mathbf{x}_0)}{\delta}$$

Une valeur importante de $|S_j^{\mathbf{x}_0}|$ est révélateur d'un modèle sensible aux variations de l'entrée X_j .

Inconvénient : indice local, dépend du choix de \mathbf{x}_0 .

- Méthode de Morris. Indices de sensibilité globaux :

$$\mu_j = \mathbb{E}[|S_j^{\mathbf{X}}|], \quad \sigma_j^2 = \text{Var}(S_j^{\mathbf{X}})$$

Estimation des indices :

- On tire R points $(\mathbf{x}_k)_{k=1,\dots,R}$ (avec Monte Carlo, ou bien avec quasi-Monte Carlo; par exemple, avec des X_i uniformes sur $[0, 1]$ et indépendants, on tire les \mathbf{x}_k avec un hypercube latin).

- On fait $R(d + 1)$ appels au code en \mathbf{x}_k et $\mathbf{x}_k + \delta \mathbf{e}_j$, $j = 1, \dots, d$, $k = 1, \dots, R$.

$$\hat{\mu}_j = \frac{1}{R} \sum_{k=1}^R |S_j^{\mathbf{x}_k}|, \quad \hat{\sigma}_j^2 = \frac{1}{R-1} \sum_{k=1}^R \left(S_j^{\mathbf{x}_k} - \frac{1}{R} \sum_{l=1}^R S_j^{\mathbf{x}_l} \right)^2$$

$$\mu_j = \mathbb{E}[|S_j^{\mathbf{X}}|] \text{ estimé par } \hat{\mu}_j = \frac{1}{R} \sum_{k=1}^R |S_j^{\mathbf{x}_k}|$$

$$\sigma_j^2 = \text{Var}(S_j^{\mathbf{X}}) \text{ estimé par } \hat{\sigma}_j^2 = \frac{1}{R-1} \sum_{k=1}^R \left(S_j^{\mathbf{x}_k} - \frac{1}{R} \sum_{l=1}^R S_j^{\mathbf{x}_l} \right)^2$$

- μ_j est une mesure de la sensibilité.

Une valeur importante traduit des effets importants (en moyenne) et un modèle sensible aux variations de l'entrée X_j .

- σ_j est une mesure des interactions et des effets non-linéaires.

Une valeur importante traduit des effets différents les uns des autres, des effets qui dépendent de la valeur :

- soit de l'entrée elle-même : effet non-linéaire
- soit des autres entrées : interaction

(la méthode ne permet pas de distinguer les deux cas)

Criblage : Méthode de Morris

- Permet de classer les entrées en trois groupes selon leurs effets :
 - effets négligeables
 - effets linéaires et sans interaction
 - effets non-linéaires et/ou avec interactions
- En théorie, pas d'hypothèse sur le modèle; en pratique, la méthode marche mieux avec une certaine régularité.
- Généralise la méthode OAT (One factor At a Time).

Indices de sensibilité pour un modèle linéaire

- Soit le modèle linéaire

$$Y = f(X_1, \dots, X_d) = \alpha_0 + \sum_{i=1}^d \alpha_i X_i$$

avec $X_i \in \mathbb{R}$ décorrélés.

- La variance de Y s'écrit

$$\text{Var}(Y) = \sum_{i=1}^d \alpha_i^2 \text{Var}(X_i)$$

La part de variance de Y due à X_i est $\alpha_i^2 \text{Var}(X_i)$.

→ On obtient une décomposition de la variance de Y en fonction des X_i .

- La sensibilité de Y à X_i est quantifiée par l'indice SRC (Standardized Regression Coefficient) :

$$\text{SRC}_i = \frac{\alpha_i^2 \text{Var}(X_i)}{\text{Var}(Y)}$$

- Avec le modèle linéaire $Y = \alpha_0 + \sum_{i=1}^d \alpha_i X_i$,

$$\text{SRC}_i = \frac{\alpha_i^2 \text{Var}(X_i)}{\text{Var}(Y)}$$

de manière équivalente :

$$\text{SRC}_i = \frac{\text{Cov}(X_i, Y)^2}{\text{Var}(X_i) \text{Var}(Y)}$$

- Estimation :

- estimation des α_i (OAT différences finies, régression linéaire).
- Monte Carlo : on tire $\mathbf{X}^{(k)}$, on calcule $Y^{(k)} = f(\mathbf{X}^{(k)})$, et on estime

$$\widehat{\text{SRC}}_i = \frac{\left(\sum_{k=1}^n (X_i^{(k)} - \hat{X}_i)(Y^{(k)} - \hat{Y}) \right)^2}{\sum_{k=1}^n (X_i^{(k)} - \hat{X}_i)^2 \sum_{k=1}^n (Y^{(k)} - \hat{Y})^2},$$

$$\hat{X}_i = \frac{1}{n} \sum_{k=1}^n X_i^{(k)}, \quad \hat{Y} = \frac{1}{n} \sum_{k=1}^n Y_i^{(k)},$$

Indices de sensibilité pour un modèle monotone

- Soit le modèle monotone

$$Y = f(X_1, \dots, X_d)$$

avec $X_i \in \mathbb{R}$ décorrélés et f monotone.

- Coefficient de corrélation basé sur les rangs :

- On tire avec MC $\mathbf{X}^{(k)}$ et on calcule $Y^{(k)} = f(\mathbf{X}^{(k)})$, $k = 1, \dots, n$.

- A chaque k on associe son rang $r_Y^{(k)}$ selon la valeur $Y^{(k)}$ (rang 1 attribué à la plus petite valeur, rang n associé à la plus grande valeur) et son rang $r_{X_i}^{(k)}$ selon la valeur de $X_i^{(k)}$.

- La sensibilité de Y à X_i est quantifiée par l'indice SRRC (Standardized Regression Rank Coefficient) qui est le SRC pour les rangs :

$$\widehat{\text{SRRC}}_i = \frac{\left(\sum_{k=1}^n (r_{X_i}^{(k)} - \frac{n+1}{2})(r_Y^{(k)} - \frac{n+1}{2}) \right)^2}{\sum_{k=1}^n (r_{X_i}^{(k)} - \frac{n+1}{2})^2 \sum_{k=1}^n (r_Y^{(k)} - \frac{n+1}{2})^2}$$

Indices de sensibilité pour un modèle arbitraire (Sobol)

- Lorsqu'il n'est pas possible de faire d'hypothèse sur le modèle (cas général), on définit des indices de sensibilité à partir d'une décomposition de la variance de Y .
- Soit le modèle $Y = f(X_1, \dots, X_d)$ avec $X_i \in \mathbb{R}$ et $X_i \perp X_j$, $i \neq j$.
- Le modèle peut se décomposer (Sobol) en

$$f(X_1, \dots, X_d) = f_0 + \sum_{i=1}^d f_i(X_i) + \sum_{1 \leq i < j \leq d} f_{ij}(X_i, X_j) + \dots + f_{1\dots d}(X_1, \dots, X_d)$$

avec

$$\mathbb{E}[f_{i_1 \dots i_s}(X_{i_1}, \dots, X_{i_s}) f_{j_1 \dots j_t}(X_{j_1}, \dots, X_{j_t})] = 0 \text{ si } (i_1 \dots i_s) \neq (j_1 \dots j_t)$$

En fait, un seul choix possible :

$$\begin{aligned} f_0 &= \mathbb{E}[Y] \\ f_i(X_i) &= \mathbb{E}[Y|X_i] - f_0 \\ f_{ij}(X_i, X_j) &= \mathbb{E}[Y|X_i, X_j] - f_i(X_i) - f_j(X_j) - f_0 \\ &\vdots \end{aligned}$$

En utilisant cette décomposition, la variance de Y se décompose en :

$$\text{Var}(Y) = D = \sum_{i=1}^d D_i + \sum_{1 \leq i < j \leq p} D_{ij} + \cdots + D_{1\dots p}$$

où

$$D_i = \text{Var}(E[Y|X_i])$$

$$D_{ij} = \text{Var}(E[Y|X_i, X_j] - E[Y|X_i] - E[Y|X_j])$$

- D_i est la sensibilité au 1er ordre de Y par rapport à la variable X_i (la part de la variance de Y expliquée par les fluctuations de X_i)
- D_{ij} est la sensibilité de Y par rapport à l'interaction des variables X_i et X_j (la part de la variance Y expliquée par l'interaction de X_i et X_j) .
- Indices de sensibilité :

$$S_i = \frac{D_i}{D}, \quad S_{ij} = \frac{D_{ij}}{D}$$

avec

$$S_i \geq 0, \quad S_{ij} \geq 0, \quad \text{et} \quad \sum_{i=1}^d S_i + \sum_{1 \leq i < j \leq p} S_{ij} + \cdots + S_{1\dots p} = 1$$

Le nombre d'indices de sensibilité $2^d - 1$ devient vite grand avec d .

On introduit alors l'indice de sensibilité total

$$S_{T_i} = \text{somme de tous les indices relatifs à } X_i$$

qui exprime la sensibilité de Y à X_i sous toutes ses formes, i.e. à X_i seule et en interaction avec d'autres variables.

On donne en général l'indice du premier ordre S_i et l'indice total S_{T_i} .

- Si S_{T_i} petit : la variable X_i a des effets négligeables,
- Si S_i grand : la variable X_i a des effets propres importants,
- Si S_i petit et S_{T_i} grand : la variable X_i a des effets importants en interaction avec d'autres variables.

- Méthodes d'estimation des indices de Sobol.
- Méthodes particulières :
 - Mc Kay
 - FAST (Fourier Amplitude Sensitivity Test)
- Méthode de Sobol
 - basée sur des estimations de Monte Carlo,
 - des améliorations existent (avec des techniques de réduction de variance) :
 - échantillonnage stratifié,
 - échantillonnage par plans space-filling (Hypercubes Latins),
 - suites à discrédance faible (suites de Sobol $LP\tau$).
- Estimations gourmandes en nombre d'appels.
 - ↔ Utilisation d'un métamodèle (en particulier, polynômes de chaos). Mais on estime alors la sensibilité du métamodèle.

- Estimations de Sobol.

Indice de premier ordre : $S_i = \frac{D_i}{D} = \frac{\text{Var}(E[Y|X_i])}{\text{Var}(Y)}$.

$f_0 = \mathbb{E}[Y]$ et $D = \text{Var}(Y)$ sont estimés classiquement par

$$f_0 \simeq \frac{1}{n} \sum_{k=1}^n f(X_1^{(k)}, \dots, X_{i-1}^{(k)}, X_i^{(k)}, X_{i+1}^{(k)}, \dots, X_d^{(k)})$$

$$D \simeq \frac{1}{n} \sum_{k=1}^n f(X_1^{(k)}, \dots, X_{i-1}^{(k)}, X_i^{(k)}, X_{i+1}^{(k)}, \dots, X_d^{(k)})^2 - f_0^2$$

où $(X_1^{(k)}, \dots, X_d^{(k)})_{k=1 \dots n}$ est un n -échantillon des variables d'entrée.

$D_i = \text{Var}(E[Y|X_i])$ est estimée par

$$D_i \simeq \frac{1}{n} \sum_{k=1}^n f(X_1^{(k)}, \dots, X_{i-1}^{(k)}, \mathbf{X}_i^{(k)}, X_{i+1}^{(k)}, \dots, X_d^{(k)}) f(\tilde{X}_1^{(k)}, \dots, \tilde{X}_{i-1}^{(k)}, \mathbf{X}_i^{(k)}, \tilde{X}_{i+1}^{(k)}, \dots, \tilde{X}_d^{(k)}) - f_0^2$$

où $(\tilde{X}_1^{(k)}, \dots, \tilde{X}_d^{(k)})_{k=1 \dots n}$ est un second échantillon des variables d'entrée.

Questions complémentaires (plus ou moins ouvertes)

Incertitude de modèle :

Que deviennent les indices de sensibilité estimés si le modèle change (mutation) ?

↳ quelques résultats sur quelques modèles de mutation (additifs).

Comment prendre en compte, dans les résultats de sensibilité, l'utilisation d'un modèle simplifié ?

↳ pour l'estimation Monte Carlo des indices, pour des approches multi-fidélité.

Modèles à entrées corrélées :

Comment réaliser (interpréter) une analyse de sensibilité lorsque les variables d'entrée ne sont pas indépendantes ?

↳ se ramener à des variables décorréées.

Références :

A. Saltelli, K. Chan et E. M. Scott, Sensitivity analysis

A. Saltelli et S. Tarantola, Sensivity analysis in practice

Partie IV. Plans d'expériences et métamodèles

Planification d'expériences : Mise au point d'une suite d'expériences (de simulations numériques) en fonction d'un but :

- criblage (détermination des paramètres d'entrée importants)
- étude quantitative de l'influence des paramètres d'entrée
- construction d'une surface de réponse, optimisation, ...

Référence :

J. J. Dreesbeke, J. Fine et G. Saporta, Plans d'expériences, application à l'entreprise.

Ici : construction d'un métamodèle (ou surface de réponse).

Type de métamodèles (fonction f_r destinée à approcher la fonction f) :

- Polynômes
- Splines (fonctions définies par morceaux par des polynômes)
- Modèles linéaires généralisés
- Polynômes de chaos
- Krigeage (la fonction f est représentée comme une réalisation d'un processus gaussien; utilise la théorie des processus gaussiens)
- Réseaux de neurones (optimisés par méthodes d'apprentissage à des fins de mémorisation et de généralisation).

mémorisation : le fait d'assimiler des exemples éventuellement nombreux,

généralisation : le fait d'être capable, grâce aux exemples appris, de traiter des exemples distincts, encore non rencontrés, mais similaires.

- Machines à vecteurs de support (en anglais Support Vector Machine SVM) sont destinées à résoudre des problèmes de discrimination (à l'origine, classifieur linéaire).

↪ Chaque méthode a un (des) plan(s) d'expériences approprié(s).

Plans pour des surfaces de réponse polynomiales

- Supposons \mathbf{X} de dimension d à moyenne nulle. La "vraie" fonction est :

$$Y = f(\mathbf{X}) = \beta_0 + \sum_{j=1}^d \beta_j X_j + \sum_{1 \leq i < j \leq d} \beta_{ij} X_i X_j + \dots$$

$\beta_j X_j$: effet principal; $\beta_{ij} X_i X_j$: interaction d'ordre 2; etc.

Modèle polynomial (degré 1) :

$$f_r(\mathbf{X}) = b_0 + \sum_{j=1}^d b_j X_j$$

Modèle polynomial (degré 2) :

$$f_r(\mathbf{X}) = b_0 + \sum_{j=1}^d b_j X_j + \sum_{1 \leq i < j \leq d} b_{ij} X_i X_j$$

- Résolution : une méthode est dite de résolution L si aucune interaction d'ordre m n'est confondue avec des interactions d'ordre $L - m - 1$.

Résolution III : aucun effet principal (d'ordre 1) n'est confondu (aliasing) avec un autre effet principal (mais un effet principal peut être confondu avec un effet d'interaction d'ordre 2).

Plackett-Burman design $2^3//4$

”Vraie” fonction : $Y = \beta_0 + \sum_{j=1}^3 \beta_j X_j + \sum_{1 \leq i < j \leq 3} \beta_{ij} X_i X_j + \dots$

Modèle postulé : $f_r(\mathbf{X}) = b_0 + \sum_{j=1}^3 b_j X_j = (1, \mathbf{X}^T) \mathbf{b}$

Plan :

n	X_0	X_1	X_2	X_3
1	+1	-1	-1	+1
2	+1	+1	-1	-1
3	+1	-1	+1	-1
4	+1	+1	+1	+1

Le plan forme une matrice \mathbf{M} , de taille 4×4 .

On effectue les 4 simulations $\rightarrow \mathbf{Y}_{obs}$, on résoud $\mathbf{Y}_{obs} \simeq \mathbf{M}\mathbf{b}$, qui donne $\mathbf{b} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T \mathbf{Y}_{obs}$.

$$b_0 = \beta_0 + \dots$$

$$b_1 = \beta_1 + \beta_{23} + \dots$$

$$b_2 = \beta_2 + \beta_{13} + \dots$$

$$b_3 = \beta_3 + \beta_{12} + \dots$$

Résolution R_{III} .

- Résolution : une méthode est dite de résolution L si aucune interaction d'ordre m n'est confondue avec des interactions d'ordre $L - m - 1$.

Résolution III : aucun effet principal (d'ordre 1) n'est confondu (aliasing) avec un autre effet principal (mais un effet principal peut être confondu avec un effet d'interaction d'ordre 2).

Résolution IV : aucun effet principal n'est confondu avec un autre effet principal ou un effet d'ordre 2 (mais un effet principal peut être confondu avec un effet d'interaction d'ordre 3, ou bien un effet d'ordre 2 peut être confondu avec un autre effet d'ordre 2).

Résolution V : aucun effet principal n'est confondu avec un autre effet principal ou un effet d'ordre 2 ou 3, aucun effet d'ordre 2 n'est confondu avec un autre effet principal ou un effet d'ordre 2.

- Plackett-Burman (d paramètres d'entrée à deux niveaux)
 - $n \geq d + 1$, $n \equiv 0[4]$, R_{III}
 - $n \geq 2(d + 1)$, $n \equiv 0[4]$, R_{IV}
- Plackett-Burman (q_j niveaux pour le paramètre j , $j = 1, \dots, d$)
 - $n \geq 1 + \sum_{j=1}^d (q_j - 1)$, R_{III}
- Design de Rechtschaffner (d paramètres d'entrée à deux niveaux).
 - $n = 1 + d(d - 1)/2$, $R_V \rightarrow$ évalue tous les effets principaux et toutes les interactions d'ordre 2

Plans pour des surfaces de réponse polynomiales

- Pas de plan d'expérience universel !

Dépend de la forme de la surface de réponse et du domaine de variation (loi) des X_i .

- Plusieurs sens d'optimalité sont possibles (souvent : minimiser la variance de prédiction intégrée sur le domaine).
- Nécessité de valider sur des points de validation (ou de test), et donc de construire la surface de réponse avec seulement une partie des simulations/des expériences à disposition.

Plans space-filling

Quand la forme de la surface n'est pas a priori connue (en particulier, pour des surfaces de réponse non-paramétriques de type krigeage) : plans space-filling (pour balayer le domaine d'intérêt).

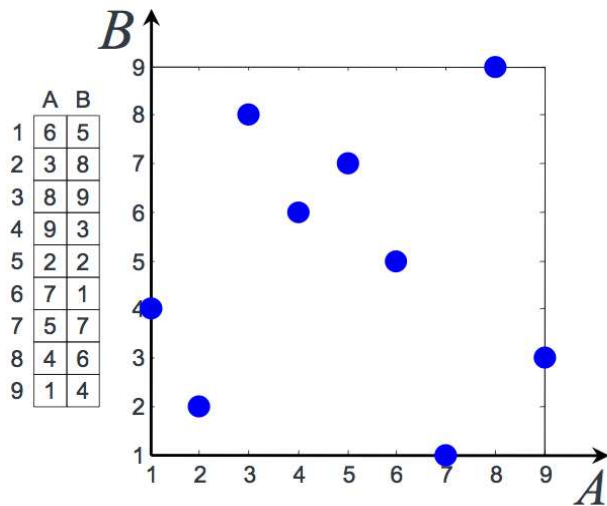
- Typiquement hypercubes latins (HL) (domaine d'intérêt : hypercube) :
 - Chaque facteur a le même nombre de niveaux (n).
 - Chacun des niveaux est pris une fois et une seule par chaque facteur.
- Critères de sélection parmi tous les HL $(\mathbf{x}_j)_{j=1,\dots,n}$:
 - Remplissage (maximin) : indicateur $D_{min} = \min_{1 \leq i \neq j \leq n} d(\mathbf{x}_i, \mathbf{x}_j)$ à maximiser.
 - Indépendance des facteurs : déterminant R de la matrice de corrélation à maximiser.
 - Uniformité (discrépance) : distance (centered L^2 discrepancy) à la répartition uniforme à minimiser :

$$CL_2 = \left(\frac{13}{12}\right)^d - \frac{2}{n} \sum_{i=1}^n \prod_{k=1}^d \left(1 + \frac{|x_{ik}-0.5|}{2} - \frac{|x_{ik}-0.5|}{2}\right) + \frac{1}{n^2} \sum_{i,j=1}^n \prod_{k=1}^d \left(1 + \frac{|x_{ik}-0.5|}{2} + \frac{|x_{jk}-0.5|}{2} - \frac{|x_{ik}-x_{jk}|}{2}\right)$$

- Méthodes de sélection :

- Méthode exploratoire : On génère un grand nombre de plans et on retient le meilleur selon le critère retenu.

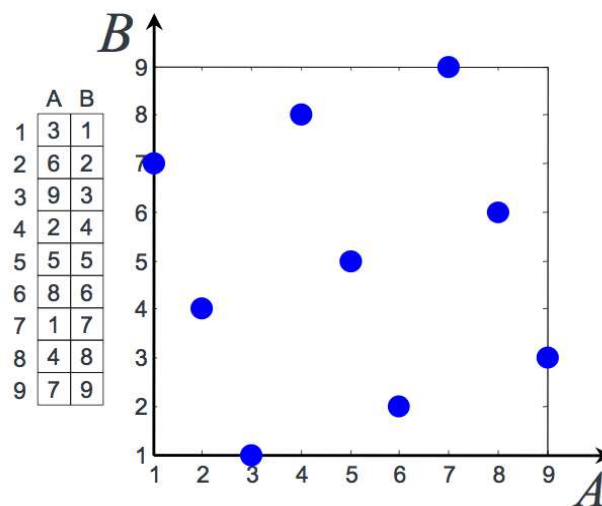
- Méthodes évoluées : Il existe dans la littérature différentes heuristiques (recuit simulé, algorithmes d'échanges) pour optimiser les plans selon le critère que l'on souhaite.



$$D_{min} = 1.4142$$

$$R = 0.9989$$

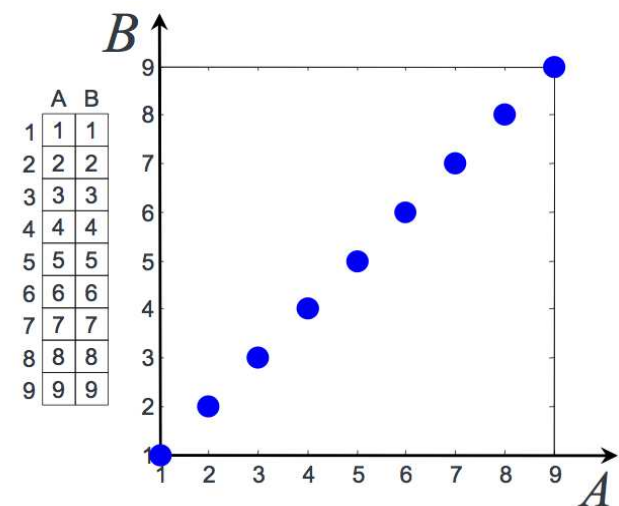
$$CL_2 = 4.19e-3$$



$$D_{min} = 3.1623$$

$$R = 1$$

$$CL_2 = 3.79e-3$$



$$D_{min} = 1.4142$$

$$R = 0$$

$$CL_2 = 13.68e-3$$

↪ on recommande souvent le critère maximin.

Polynômes de chaos

- Représentation de la solution $Y = f(\mathbf{X})$ sous la forme d'une série $\sum_j a_j \Phi_j(\mathbf{X})$ construite grâce à une projection sur une base spectrale; les coefficients a_j sont calculés par le biais de projections.

R. Ghanem et P. Spanos, Stochastic finite elements: a spectral approach, Dover Publications, 2004.

- Polynôme de chaos (de Wiener).

Supposons que $\mathbf{X} = (X_i)_{i=1,\dots,d}$ avec X_i variables aléatoires gaussiennes centrées réduites indépendantes.

Les polynômes d'Hermite (tensorisés) $\phi_{\alpha}(\mathbf{x}) = \phi_{\alpha_1}(x_1) \cdots \phi_{\alpha_d}(x_d)$ sont orthonormaux dans $L^2(\mathbb{P})$:

$$\mathbb{E}[\phi_{\alpha}(\mathbf{X})\phi_{\alpha'}(\mathbf{X})] = \begin{cases} 1 & \text{si } \alpha = \alpha' \\ 0 & \text{sinon} \end{cases}$$

La variable aléatoire $Y = f(\mathbf{X})$ (sous réserve que $\mathbb{E}[Y^2] < \infty$) s'écrit

$$Y = \sum_{\alpha \in \mathbb{N}^d} a_{\alpha} \phi_{\alpha}(\mathbf{X})$$

→ La variable aléatoire Y est caractérisée par les constantes a_{α} .

En renumérotant les polynômes $(\Phi_j)_{j \in \mathbb{N}} = (\phi_{\alpha_1, \dots, \alpha_d})_{\alpha_1, \dots, \alpha_d \in \mathbb{N}}$:

$$Y = \sum_{j=0}^{\infty} a_j \Phi_j(\mathbf{X}) \simeq \sum_{j=0}^M a_j \Phi_j(\mathbf{X})$$

d : dimension de l'espace probabiliste

q : ordre le plus élevé du polynôme pour représenter Y

$M + 1 = (d + q)!/d!q!$ termes dans la somme

- Estimation des coefficients des polynômes de chaos.

$$Y = f(\mathbf{X})$$

1/ On suppose que les entrées aléatoires sont des variables aléatoires (gaussiennes) indépendantes $(X_i)_{i=1}^d$ (ou bien on se ramène à ce cas).

2/ On écrit la solution sous forme d'une somme finie de PC :

$$Y = \sum_{i=0}^M a_i \Phi_i(\mathbf{X})$$

on substitue dans l'équation, et on projette sur la base des polynômes orthogonaux :

$$\mathbb{E} \left[\Phi_k(\mathbf{X}) \sum_{i=0}^M a_i \Phi_i(\mathbf{X}) \right] = \mathbb{E} [\Phi_k(\mathbf{X}) f(\mathbf{X})] , \quad k = 0, \dots, M$$

$$a_k = \mathbb{E} [\Phi_k(\mathbf{X}) f(\mathbf{X})] , \quad k = 0, \dots, M$$

→ estimation des a_k par régression ou par quadratures numériques (déterministes ou Monte Carlo ou quasi Monte Carlo).

- Estimation des coefficients des polynômes de chaos par moindres carrés.

$Y = f(\mathbf{X})$ à approcher par $\mathbf{a}^T \Phi(\mathbf{X})$, avec $\Phi(\mathbf{x}) = (\Phi_i(\mathbf{x}))_{i=0,\dots,M}$, $\mathbf{a} = (a_i)_{i=0,\dots,M}$

Idée : minimisation de l'erreur résiduelle :

$$\hat{\mathbf{a}}_{LS} = \underset{\mathbf{a} \in \mathbb{R}^{M+1}}{\operatorname{argmin}} \mathbb{E}[(\mathbf{a}^T \Phi(\mathbf{X}) - f(\mathbf{X}))^2]$$

Si on a un n -échantillon $(\mathbf{X}^{(k)})_{k=1,\dots,n}$, estimation par

$$\hat{\mathbf{a}} = \underset{\mathbf{a} \in \mathbb{R}^{M+1}}{\operatorname{argmin}} \frac{1}{n} \sum_{k=1}^n (\mathbf{a}^T \Phi(\mathbf{X}^{(k)}) - f(\mathbf{X}^{(k)}))^2$$

c'est-à-dire :

$$\hat{\mathbf{a}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{F}, \quad A_{ki} = \Phi_i(\mathbf{X}^{(k)}), \quad F_k = f(\mathbf{X}^{(k)})$$

- Estimation de l'erreur.

L'erreur de généralisation du modèle $\hat{\mathbf{a}}^T \Phi(\mathbf{x})$ obtenu avec $(\mathbf{X}^{(k)})_{k=1, \dots, n}$ est :

$$\mathcal{E}_{gen} = \int_{\mathbb{R}^d} (\hat{\mathbf{a}}^T \Phi(\mathbf{x}) - f(\mathbf{x}))^2 \mathbb{P}(d\mathbf{x})$$

⚠ Elle est très mal estimée (risque d'overfitting) par l'erreur empirique utilisant l'ensemble d'apprentissage :

$$\mathcal{E}_{emp} = \frac{1}{n} \sum_{k=1}^n (\hat{\mathbf{a}}^T \Phi(\mathbf{X}^{(k)}) - f(\mathbf{X}^{(k)}))^2$$

↪ Validation avec l'erreur empirique sur un ensemble test, indépendant de l'ensemble d'apprentissage (bien, mais cher) :

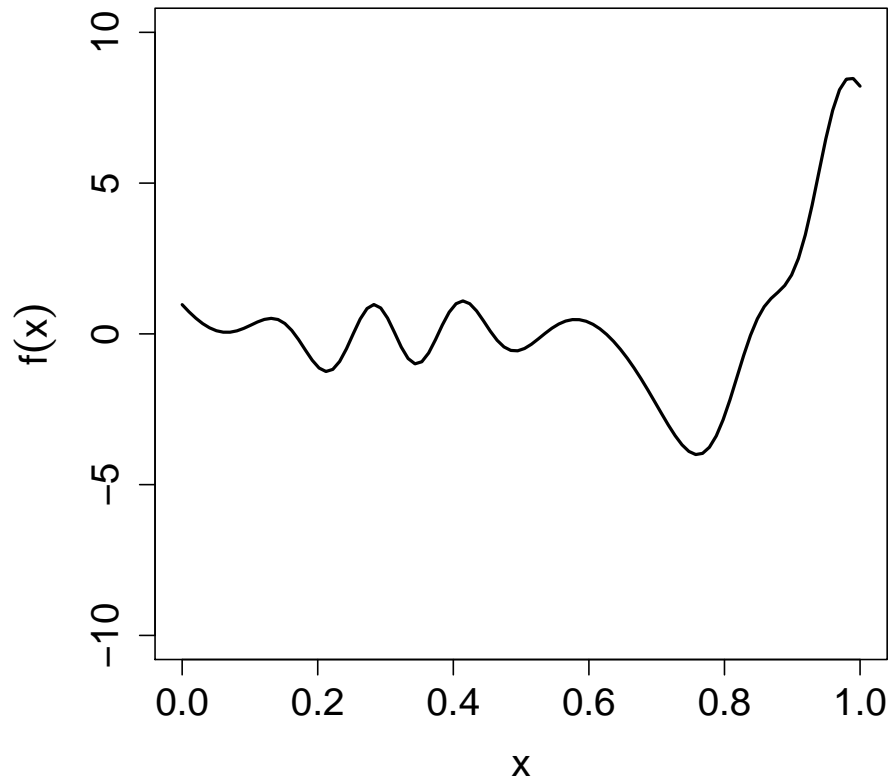
$$\mathcal{E}_{pre} = \frac{1}{n} \sum_{k=1}^n (\hat{\mathbf{a}}^T \Phi(\widetilde{\mathbf{X}}^{(k)}) - f(\widetilde{\mathbf{X}}^{(k)}))^2$$

↪ Validation par validation croisée (Leave-One Out) (moins bien, mais gratuit) :

$$\begin{aligned} \mathcal{E}_{LOO} &= \frac{1}{n} \sum_{k=1}^n \left(\hat{\mathbf{a}}[\mathbf{X}^{(-k)}]^T \Phi(\mathbf{X}^{(k)}) - f(\mathbf{X}^{(k)}) \right)^2 \\ &= \frac{1}{n} \sum_{k=1}^n \left(\frac{\hat{\mathbf{a}}^T \Phi(\mathbf{X}^{(k)}) - f(\mathbf{X}^{(k)})}{1 - q^{(k)}} \right)^2 \text{ avec } q^{(k)} = (\mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T)_{kk} \end{aligned}$$

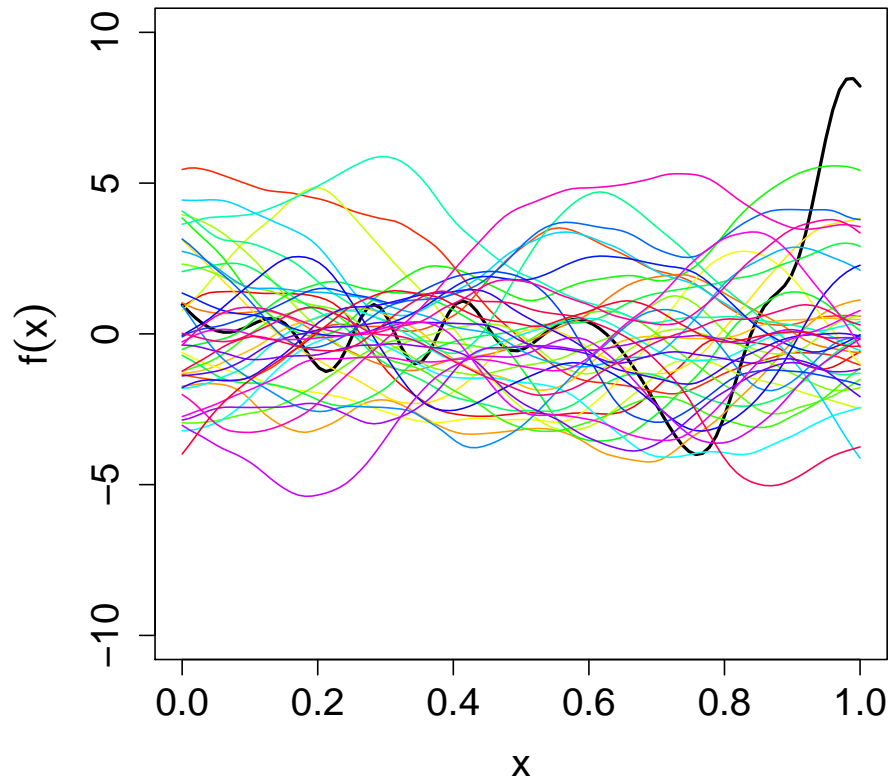
- Quelques remarques.
- Décomposition bien adaptée pour représenter les entrées ou une combinaison linéaire des entrées, mais pas forcément une fonction non-linéaire de ceux-ci. La décomposition converge d'autant plus vite que la fonction est régulière.
- Travail consiste à estimer les coefficients $(a_j)_j$ par méthodes intrusives ou non-intrusives.
- Généralisation à des distributions autres que gaussiennes. Les polynômes ϕ . sont alors modifiés.
- Avantages :
 - Calcul efficace de la sensibilité de la solution aux paramètres d'entrée incertains
 - Obtention d'une forme explicite de la solution (calcul facile de moments et de densité de probabilité)
- Inconvénients.
 - Mauvais contrôle de l'erreur (de troncation).
 - Méthode mal adaptée à des cas où le nombre de paramètres d'entrée incertains est grand.

Krigeage



- **But** : construire un métamodèle de $f(x)$.

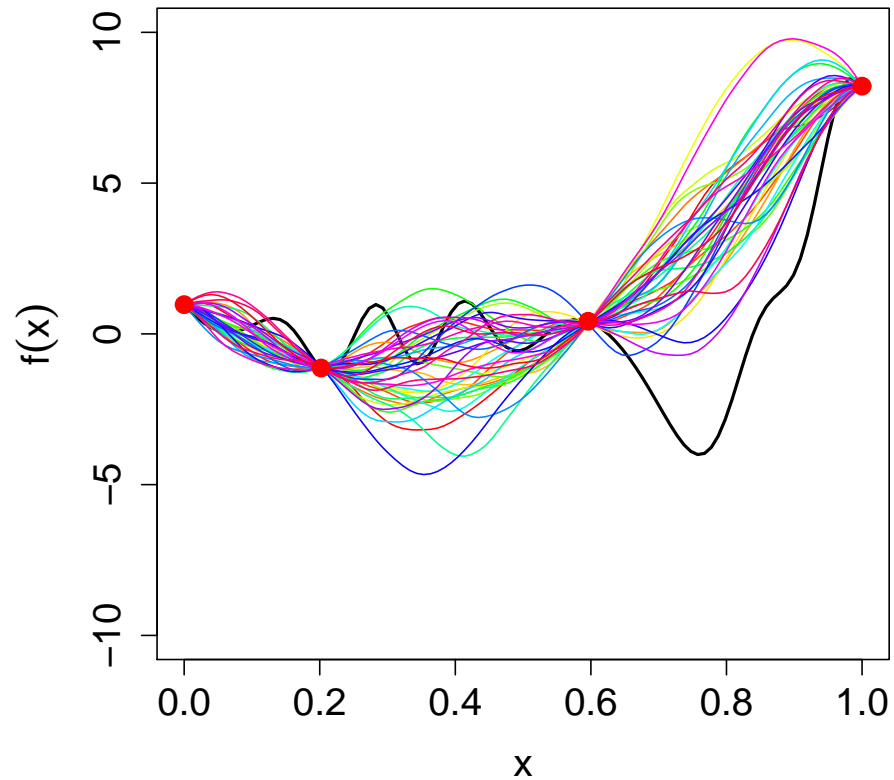
Krigeage



- **Idée** : on suppose que $f(x)$ une réalisation d'un processus gaussien $Z(x)$ avec moyenne $m(x)$ et covariance $k(x, x')$.

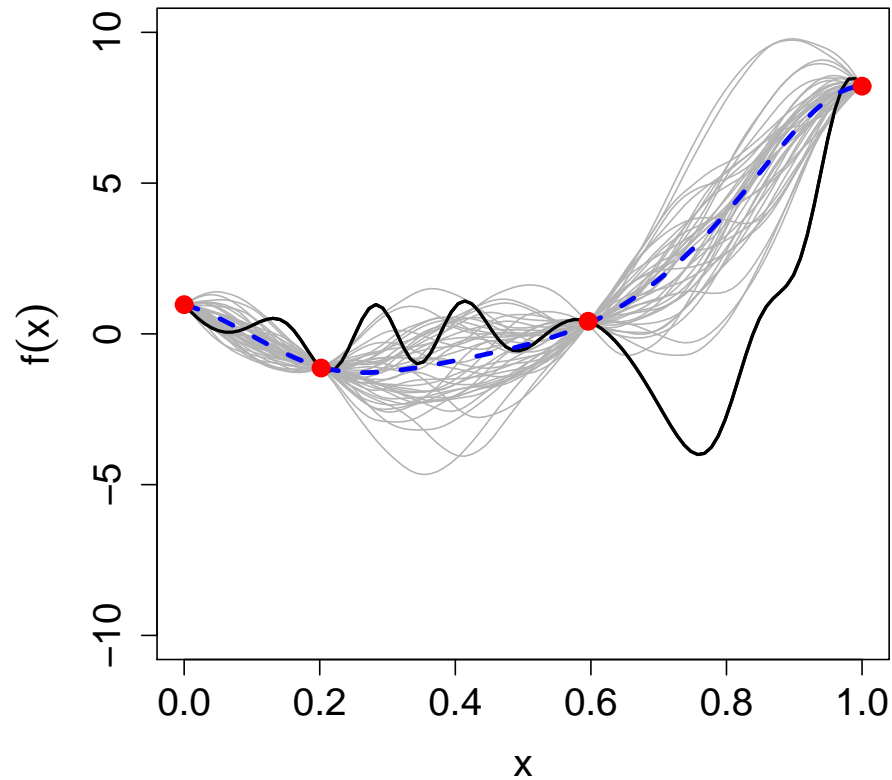
Par exemple $m(x) = \beta_0$ et $k(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{\theta^2}\right)$.

Krigeage



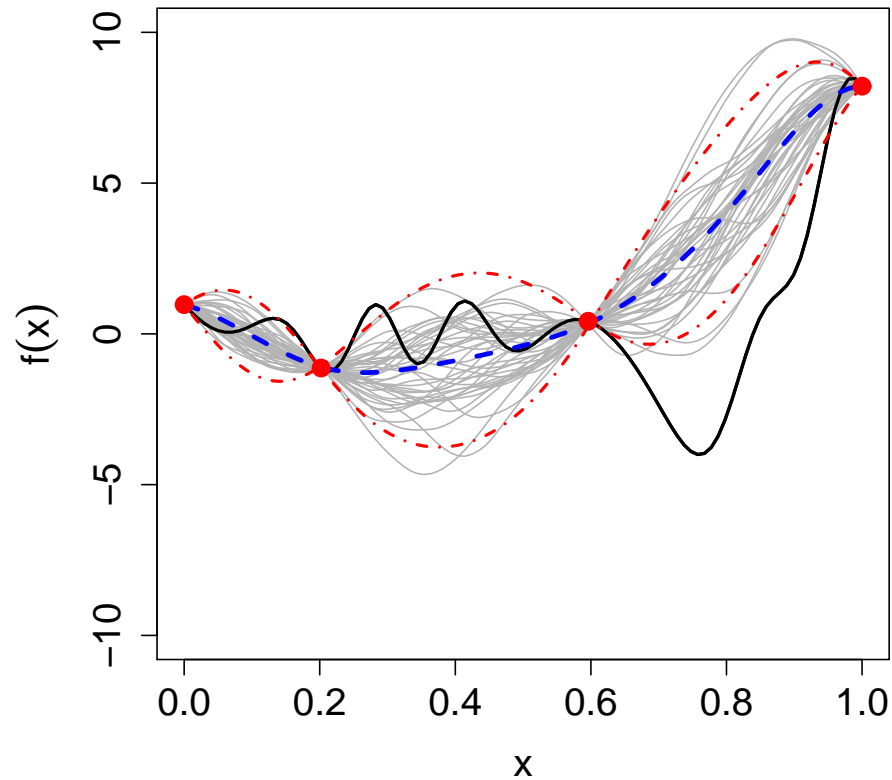
- On suppose qu'on observe $Z(x)$ aux points $x_1, \dots, x_n \in \mathbb{R}$.
- On veut prédire $Z(x)$ étant données les observations $Z(x_1), \dots, Z(x_n)$.

Krigeage



- On utilise le **Best Linear Unbiased Predictor**:
(BLUP) : $\hat{f}(x) = m(x) + \mathbf{k}(x)^T \mathbf{K}^{-1} \mathbf{z}^n$.
où $\mathbf{k}(x) = [k(x, x_i)]_{i=1, \dots, n}$,
 $\mathbf{K} = [k(x_i, x_j)]_{i, j=1, \dots, n}$,
 $\mathbf{z}^n = (f(x_1), \dots, f(x_n))$.
- **Ligne bleue : BLUP**

Krigeage



- Le BLUP est le meilleur parce qu'il minimise l'erreur quadratique moyenne

$$\hat{\sigma}^2(x) = k(x, x) - \mathbf{k}(x)^T \mathbf{K}^{-1} \mathbf{k}(x)$$

- Ligne bleue : BLUP
- Lignes rouges : intervalles de confiance

- Quelques remarques.

- Il faut déterminer la tendance $m(x)$ et la fonction de covariance $k(x, x')$.

- Pour la tendance, on considère $m(x) = \sum_{j=1}^p \beta_j f_j(x)$ avec f_j données (monômes par exemple) et on estime les β_j (avec leur incertitude).

- Pour la fonction de covariance, on suppose qu'elle appartient à une famille paramétrique $k(x, x'; \theta)$ et on estime θ par maximum de vraisemblance (méthode plug-in) ou bien on utilise une approche bayésienne complète (lourd).

- Avantages :

- Fournit des intervalles de confiance, des erreurs de généralisation.

- Obtention d'une forme explicite du BLUP.

- Permet l'implémentation de stratégie d'échantillonnage séquentielle, pour l'optimisation par exemple.

- Permet de traiter des cas bruités ou des fonctions peu régulières (la famille de fonctions de covariance peut comporter un paramètre de régularité, par exemple la famille de Matérn).

- Inconvénients :

- Méthode mal adaptée lorsque le nombre de paramètres d'entrée est grand.

- Méthode lourde lorsque le nombre de points est grand.

C. E. Rasmussen and C. Williams, Gaussian Processes for Machine Learning.