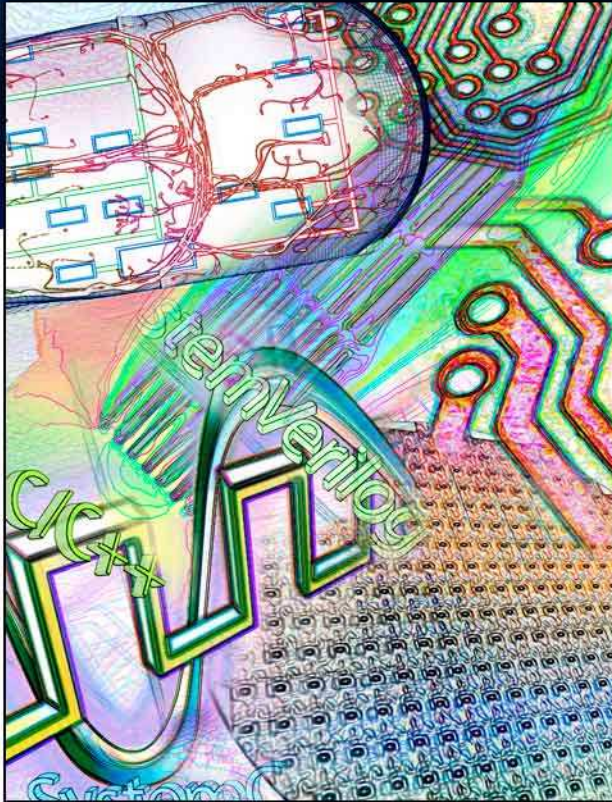


# Méthodes d'analyse de la variabilité des circuits intégrés



Hubert FILIOL

Ingénieur R&D

Mentor Graphics, Eldo-AMS

Journée thématique GDR - Méthodes et Outils pour la Prise en Compte de la Variabilité des Procédés de Fabrication

5 Octobre 2012

**Mentor**  
**Graphics**<sup>®</sup>

# Plan

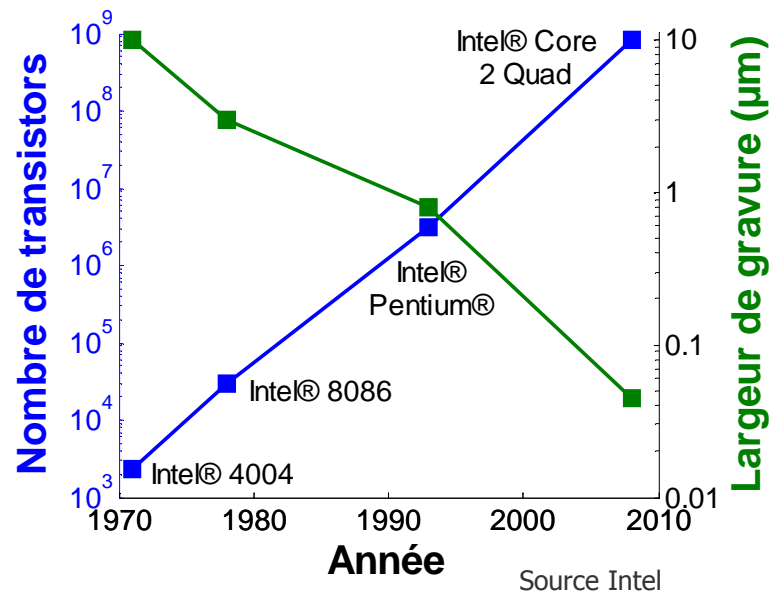
---

1. Loi de Moore et phénomènes de variabilité
2. Méthodes de caractérisation de la variabilité
3. Méthodes d'analyse de sensibilité

# LOI DE MOORE ET PHÉNOMÈNES DE VARIABILITÉ

# Evolution de la microélectronique

## ■ Loi de Moore



## ■ Intérêts

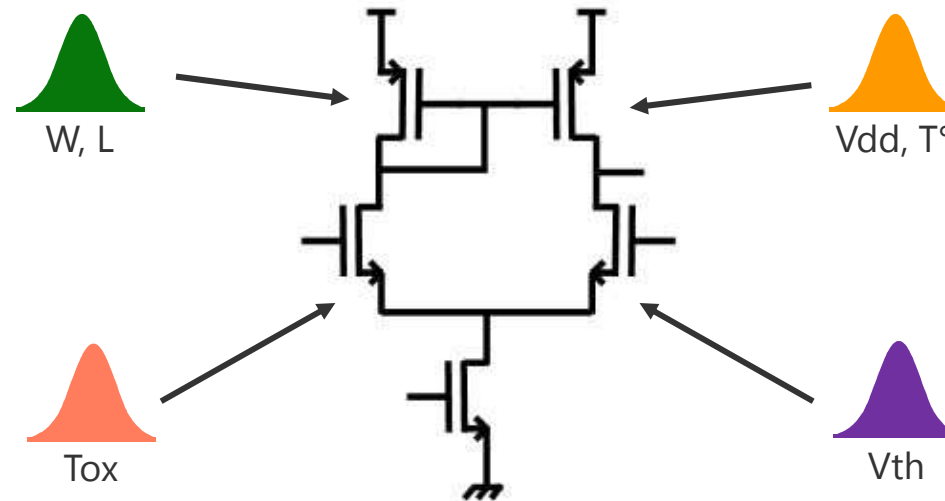
- Réduction de l'encombrement
- Nouvelles fonctionnalités
- Baisse du coût unitaire

## ■ Limites physiques

- Lithographie
- Courants de fuite
- Puissance dissipée
- Fiabilité
- ...
- **Variabilité**

# Phénomènes de variabilité

## ■ Variations des caractéristiques des circuits



## ■ Sources des variations

**Variations Paramétriques**

**Dimensions**  
**Concentration des dopants**

**Variations Environnementales**

**Tension d'alimentation**  
**Température**

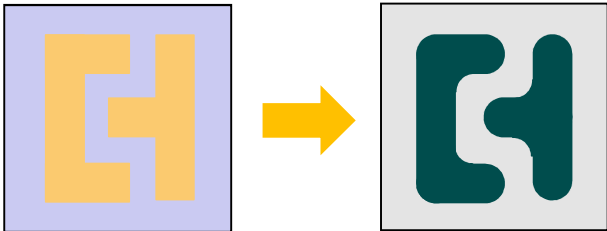
**Vieillessement**

**Electromigration**  
**Injection de porteurs chauds**

# Classification des variations

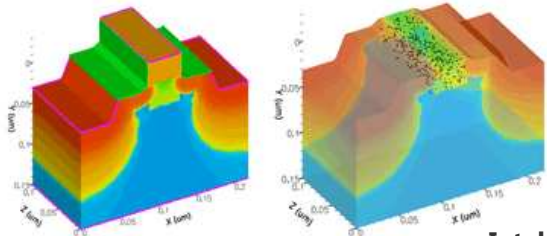
- Nature

## Variations Systématiques



Ex : lithographie, dessin des masques

## Variations Aléatoires

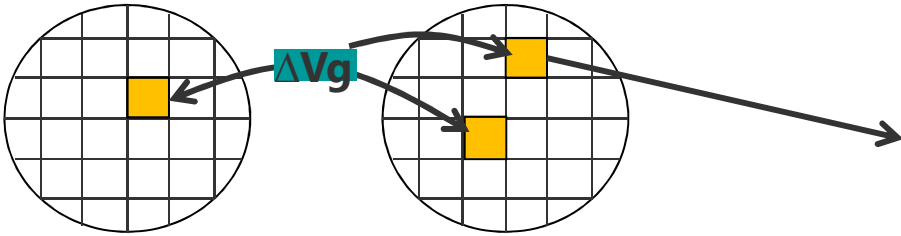


Intel

Ex : Fluctuation des dopants

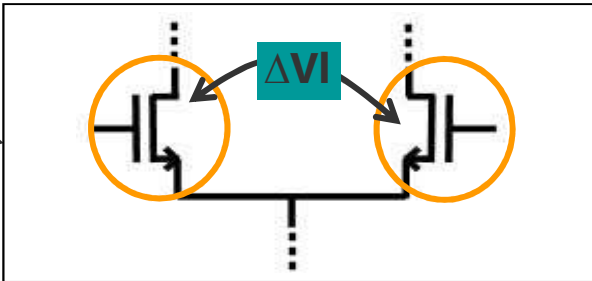
- Répartition spatiale

## Variations Globales



Inter-puces, inter-wafers, inter-lot, etc.

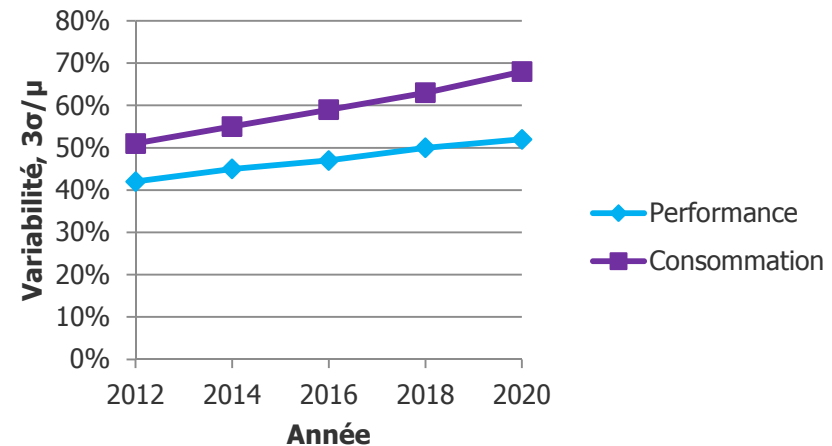
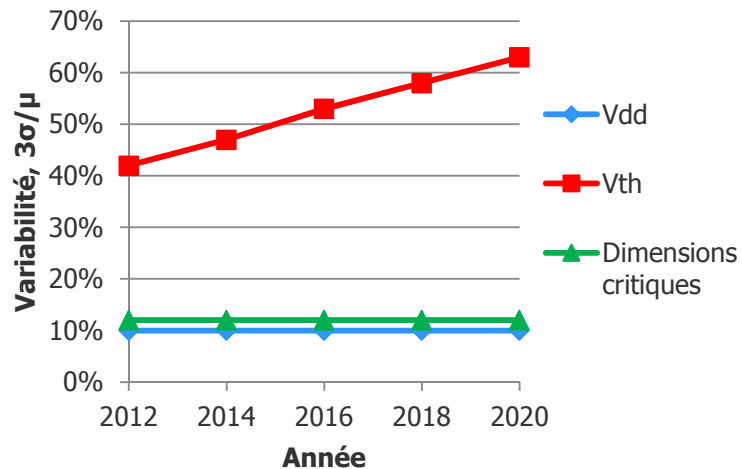
## Variations Locales



Intra-puce, mismatch

# Evolution de la variabilité

## ■ Prévision de l'ITRS (International Technology Roadmap for Semiconductors)

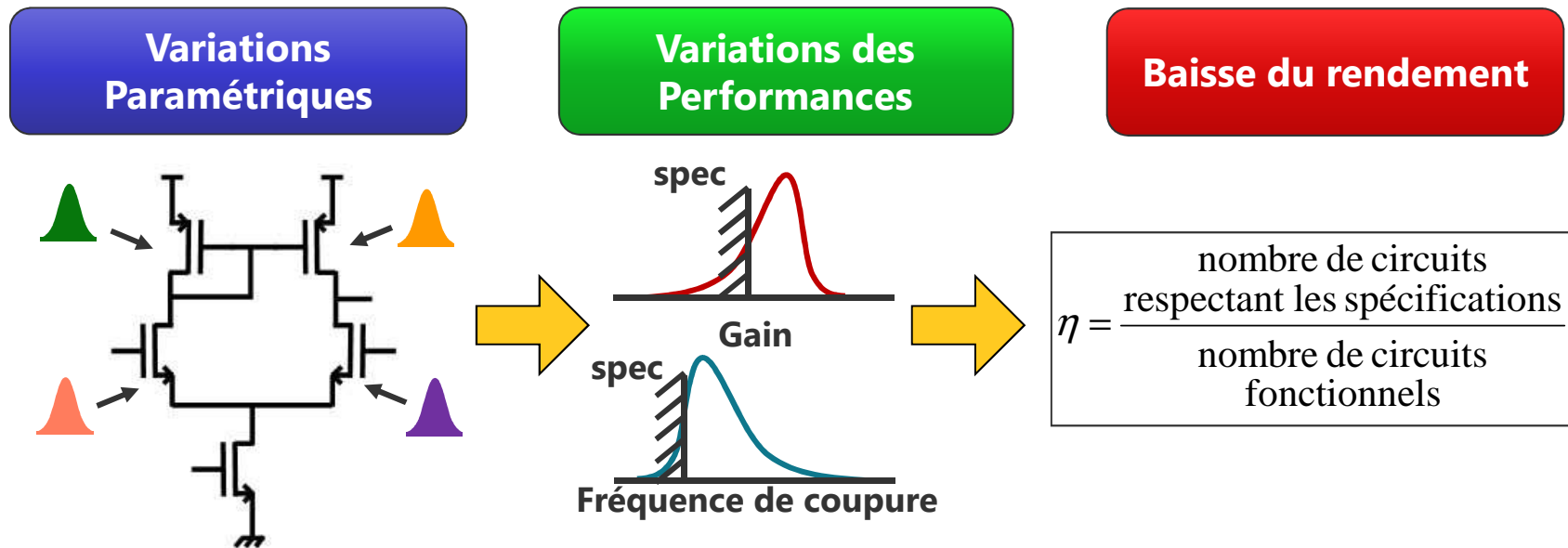


## ■ Bilan

- ↗ de la variabilité avec la réduction des dimensions
- **Variations systématiques**
  - Techniques lithographiques de compensation (OPC, PSM, etc.)
- **Variations aléatoires (RDF)**
  - Ne peuvent pas être compensées lors de la fabrication
  - Doivent être prises en compte dès le début du flot de conception

# Problématique

## ■ Baisse du rendement



## → Méthodes d'analyse de la variabilité

- Analyser l'impact des variations paramétriques sur les performances des circuits



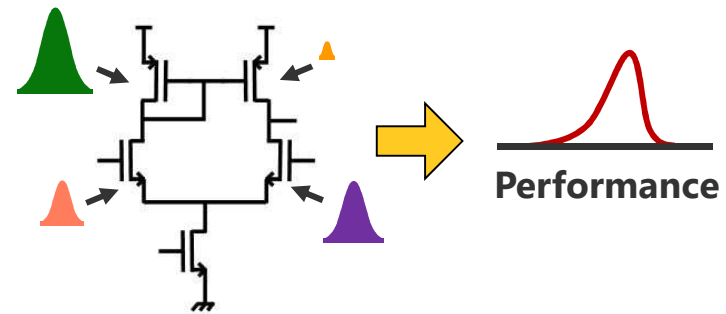
# Méthodes d'analyse de la variabilité

## Caractériser les variations des performances



- Variable aléatoire
  - Densité de probabilité
  - Moments
  - ...
- ➔ **Caractérisation de la variabilité = Estimation d'une variable aléatoire**
  - Histogramme
  - Variance

## Identifier les principales variations paramétriques

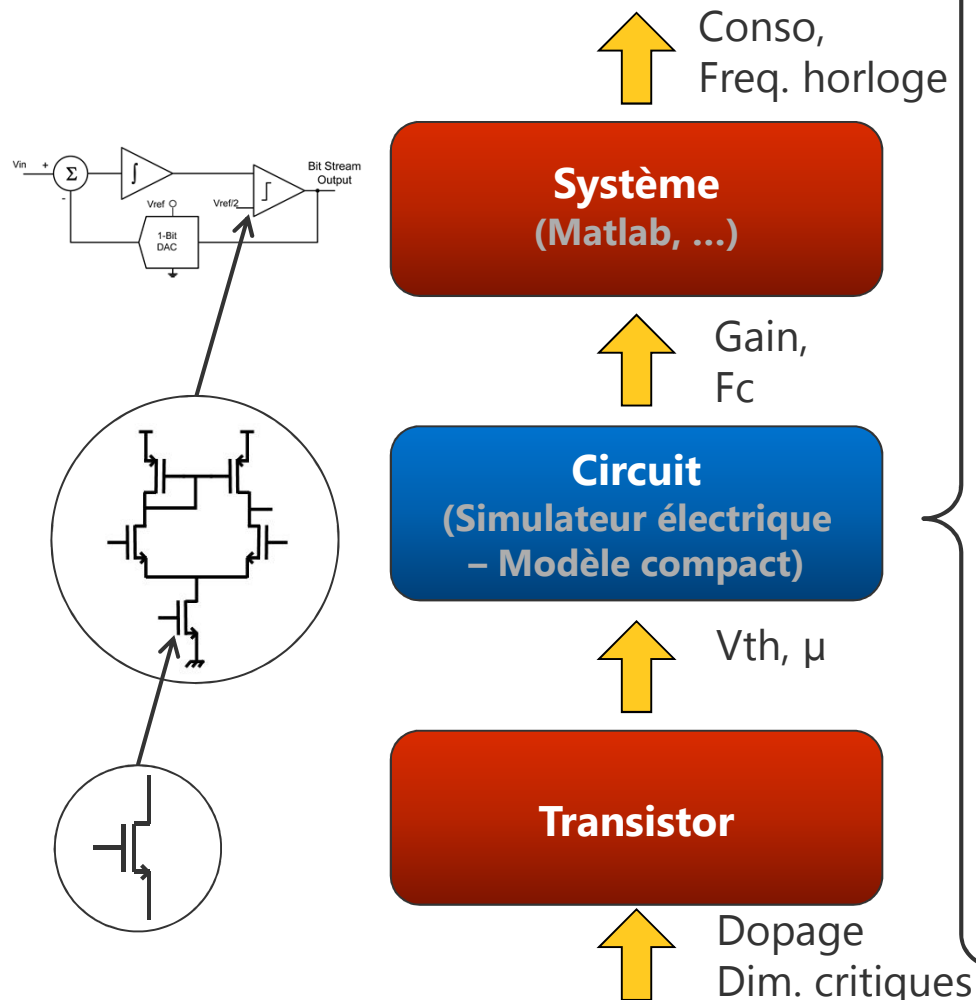


- Quelles variations impactent le plus les performances ?
- ➔ **Analyse de sensibilité**
  - Indices de sensibilité

# MÉTHODES DE CARACTÉRISATION DE LA VARIABILITÉ

# Modélisation des performances des circuits

## ■ Niveaux de modélisation



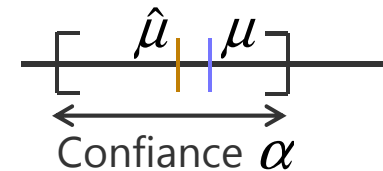
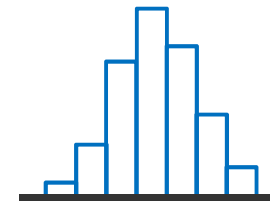
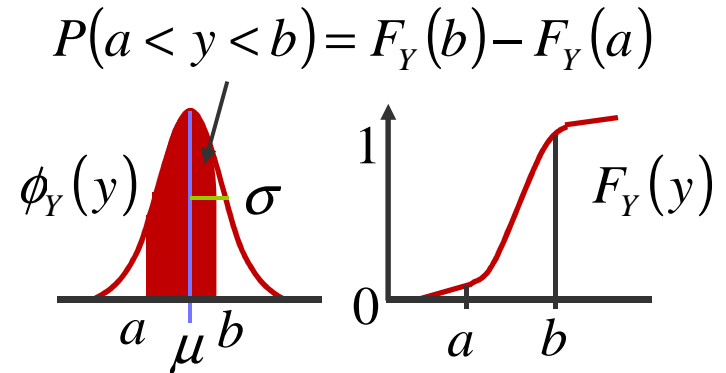
## ■ Modèle des performances :

$$Y = f(X_1, X_2, \dots, X_p)$$

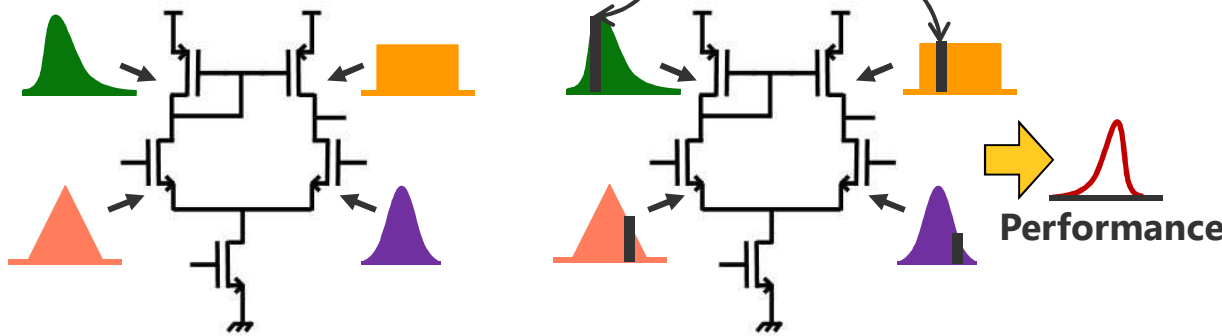
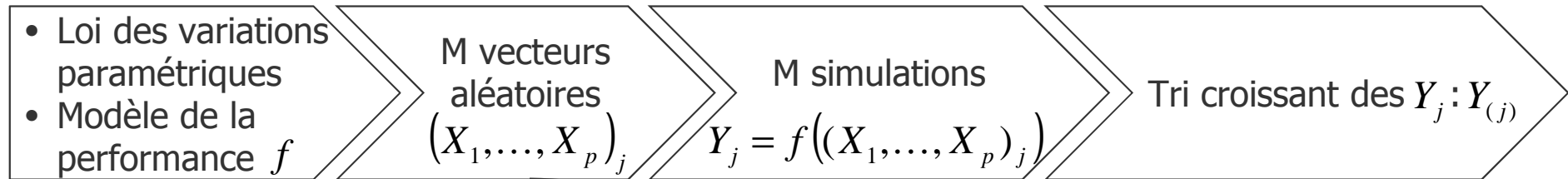
- $f$  : Simulation électrique
    - Modèle compact (BSIM, ...)
    - Spice, Eldo, ...
  - $X_i$  : Variations paramétriques
    - définies par la carte modèle d'une technologie
    - Variables aléatoires
      - gaussiennes / uniformes
      - corrélées / indépendantes
- $Y =$  Variable aléatoire

# Variables aléatoires – Estimation des variations des performances

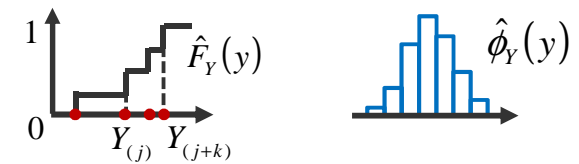
- Variable aléatoire définie par :
  - Densité de probabilité / Fonction de répartition
  - Moments (espérance, variance, ...)
- Estimation des paramètres d'une variable aléatoire + Intervalle de confiance (précision)
  - Densité de probabilité  $\Leftrightarrow$  Histogramme
  - Espérance  $\Leftrightarrow$  Moyenne arithmétique



# Méthode de Monte Carlo



- Fonction de répartition, Histo



- Moyenne, variance estimée :

$$\hat{\mu}_Y = \frac{1}{M} \sum_{j=1}^M Y_j, \quad \hat{\sigma}_Y^2 = \frac{1}{M-1} \sum_{j=1}^M (Y_j - \hat{\mu}_Y)^2$$

## Hypothèses

- Nombre de tirages M suffisamment grand

## Points +

- Générale
  - $f$  non-linéaire
  - lois asymétriques
- Précision indép. dimension

## Limites

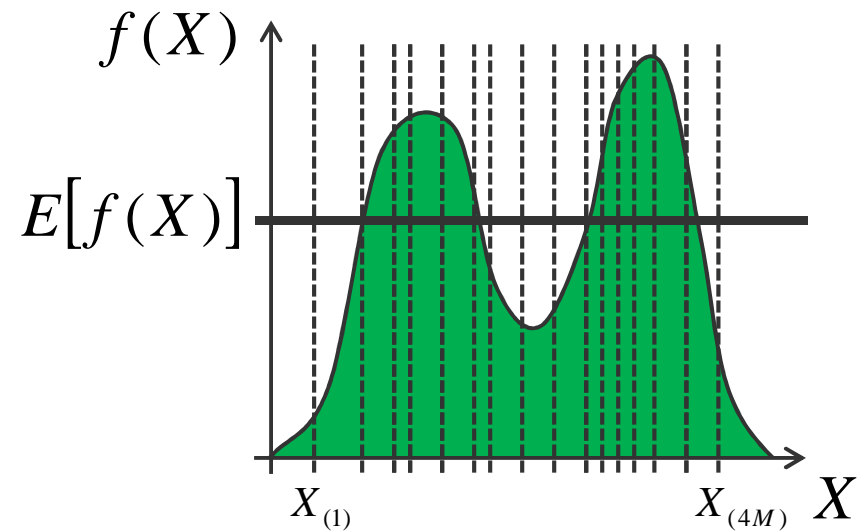
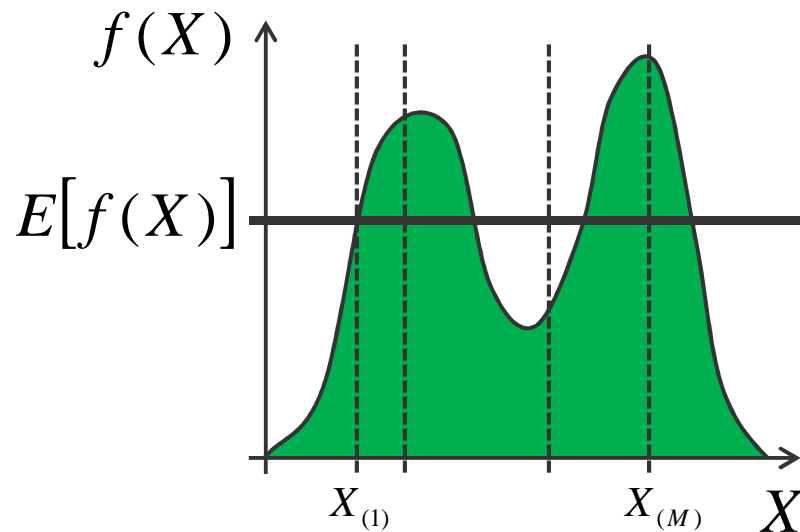
- Très lente
  - $\nearrow$  M pour  $\nearrow$  précision

# Monte Carlo – Réduction de la variance

- Soient une variable aléatoire  $X$  et une fonction  $f : \mathfrak{X} \rightarrow \mathfrak{R}$
- Objectif : calculer l'espérance  $\mu = E[f(X)]$ 
  - Générer  $M$  variables aléatoires  $(X_j)_{1 \leq j \leq M}$  indépendantes et de même loi que  $X$
  - Estimateur MC de l'espérance :  $\hat{\mu}_{MC} = \frac{1}{M} \sum_{j=1}^M f(X_j)$
- Loi des grands nombres :  $\lim_{M \rightarrow +\infty} \hat{\mu}_{MC} = \mu$
- Théorème central limite :  $M \rightarrow +\infty, \frac{\hat{\mu}_{MC} - \mu}{\sigma/\sqrt{M}} \Rightarrow N(0,1)$  avec  $\sigma^2 = \text{Var}[f(X)]$
- Estimateur de la variance :  $\hat{\sigma}^2 = \frac{1}{M-1} \sum_{j=1}^M (f(X_j) - \hat{\mu}_{MC})^2$
- Intervalle de confiance à 95% :  $\left[ \hat{\mu}_M - \frac{1.96\hat{\sigma}}{\sqrt{M}}, \hat{\mu}_M + \frac{1.96\hat{\sigma}}{\sqrt{M}} \right]$
- Améliorer la précision de l'estimation :  $\nearrow M$  ou  $\searrow \sigma$

# MC - Augmentation du nombre d'échantillons

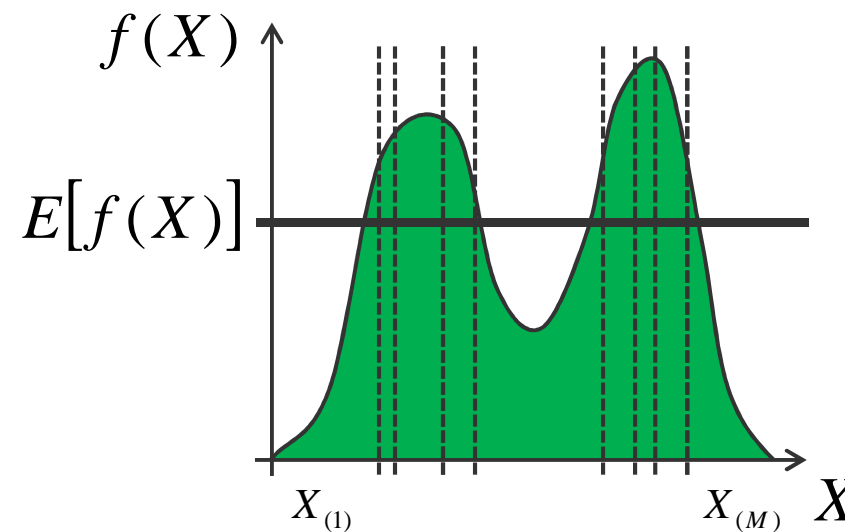
- intervalle de confiance de l'estimation proportionnel à  $1/\sqrt{M}$



- Pour diviser l'erreur par 2 il faut 4 x plus d'échantillons
- ➔ Coût calcul

# MC - Echantillonnage préférentiel (1)

- Réduction de la variance
  - Echantillonner en priorité dans les régions d'importance de la fonction





# MC - Echantillonnage préférentiel (2)

## ■ Principe : modification de l'échantillonnage

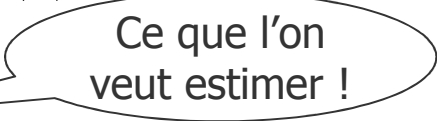
- Soit une densité de probabilité  $\tilde{\phi} > 0$
- Nouvelle écriture de l'espérance à estimer :

$$\mu = E[f(X)] = \int \frac{f(x)\phi_X(x)}{\tilde{\phi}(x)} \tilde{\phi}(x) dx = E\left[\frac{f(Z)\phi_X(Z)}{\tilde{\phi}(Z)}\right]$$

Où  $Z$  est une variable aléatoire de densité de probabilité  $\tilde{\phi}$

## ■ Variance :

$$\text{Var}\left[\frac{f(Z)\phi_X(Z)}{\tilde{\phi}(Z)}\right] = E\left[\frac{f(Z)^2\phi_X(Z)^2}{\tilde{\phi}(Z)^2}\right] - \mu^2 = \int \frac{f(x)^2\phi_X(x)^2}{\tilde{\phi}(x)} dx - \mu^2$$

→ Annulation de la variance si  $\tilde{\phi}(x) = \frac{f(x)\phi_X(x)}{\mu}$  

## ■ Heuristique :

- $\tilde{\phi}(x) \propto f(x)\phi_X(x)$

# MC - Echantillonnage préférentiel (3)

---

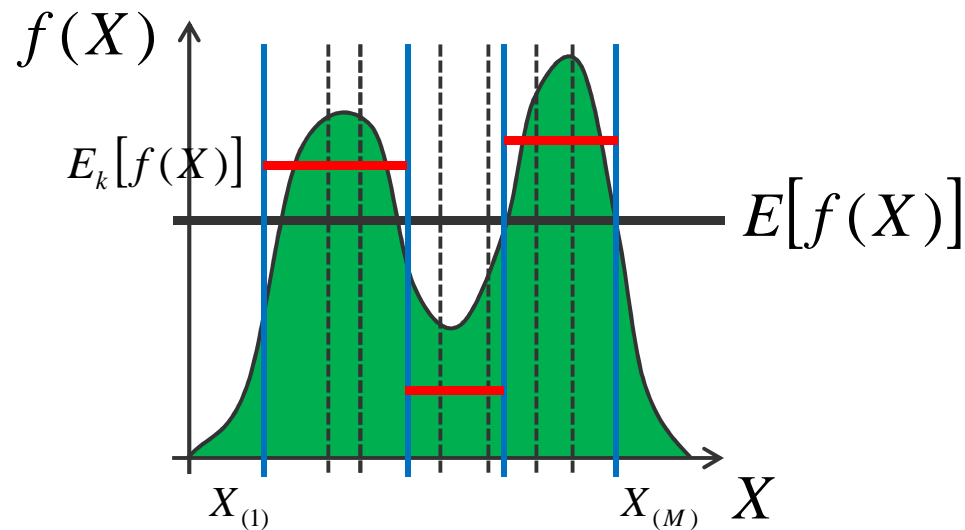
## ■ Algorithme :

- Choisir une fonction d'importance  $\tilde{\phi}$ 
  - facile à générer
  - $\tilde{\phi}(x) \propto f(x)\phi_X(x)$
  - Approche simple  $\rightarrow \tilde{\phi}$  = la distribution originale  $\phi$  mais centrée sur le maximum de  $f(x)\phi_X(x)$
- Générer  $M$  variables aléatoires  $(Z_j)_{1 \leq j \leq M}$  indépendantes de densité de probabilité  $\tilde{\phi}$
- Estimateur de l'espérance :

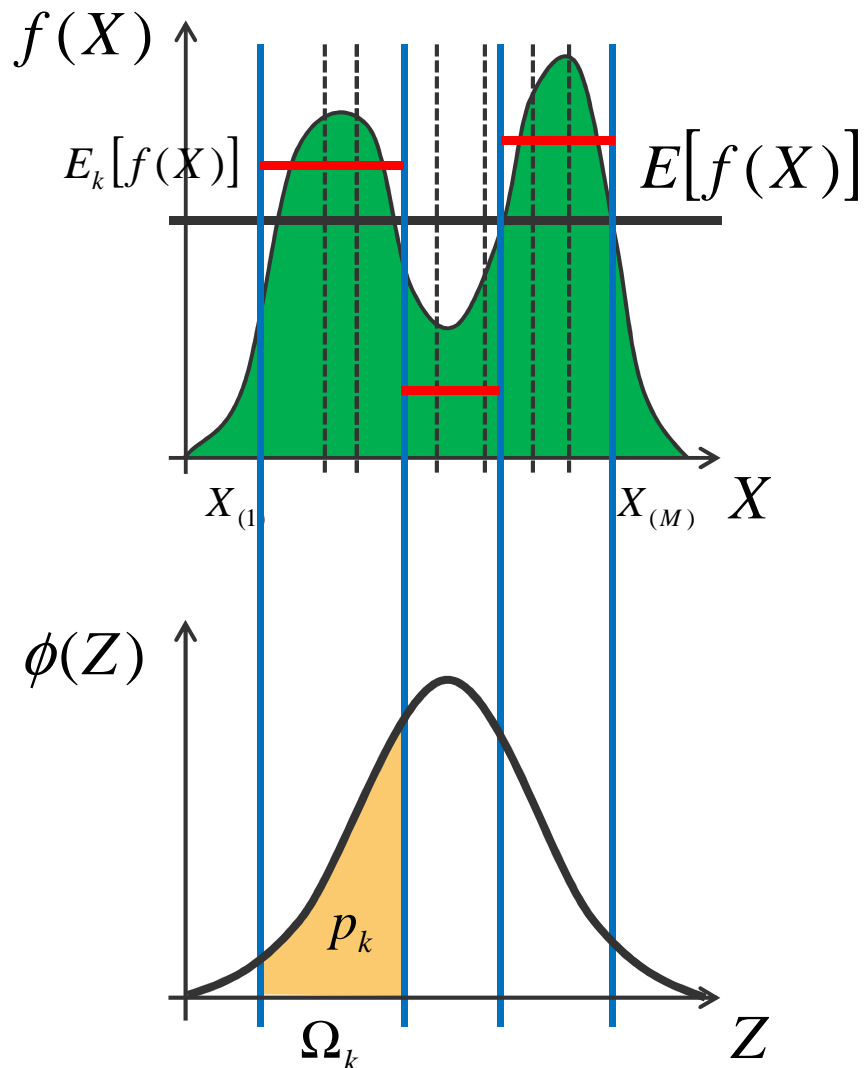
$$\hat{\mu}_{pref} = \frac{1}{M} \sum_{j=1}^M \frac{\phi_X(Z_j)}{\tilde{\phi}(Z_j)} f(Z_j)$$

# MC - Echantillonnage stratifié (1)

- Réduction de la variance :
  - Partitionner le domaine
  - Estimer l'espérance sur chaque sous-domaine avec l'estimateur MC
  - Pondérer les estimations locales par une probabilité pour obtenir l'estimation globale de l'espérance



# MC - Echantillonnage stratifié (2)



## ■ Algorithme

- Partitionner l'espace de  $Z$
- Calcul de la probabilité  $p_k$  de chaque strate
- Génération des variables aléatoires  $(X | Z \in \Omega_k)$
- Calcul de  $f(X) | Z \in \Omega_k$

## ■ Espérance

$$\mu = E[f(X)] = \sum_{k=1}^K p_k E[f(X) | Z \in \Omega_k]$$

## ■ Estimateur

$$\hat{\mu}_{strat} = \sum_{k=1}^K p_k \frac{1}{m_k} \sum_{j=1}^{m_k} f(X_j^{(k)})$$

Estimateur MC

# MC - Echantillonnage stratifié (3)

---

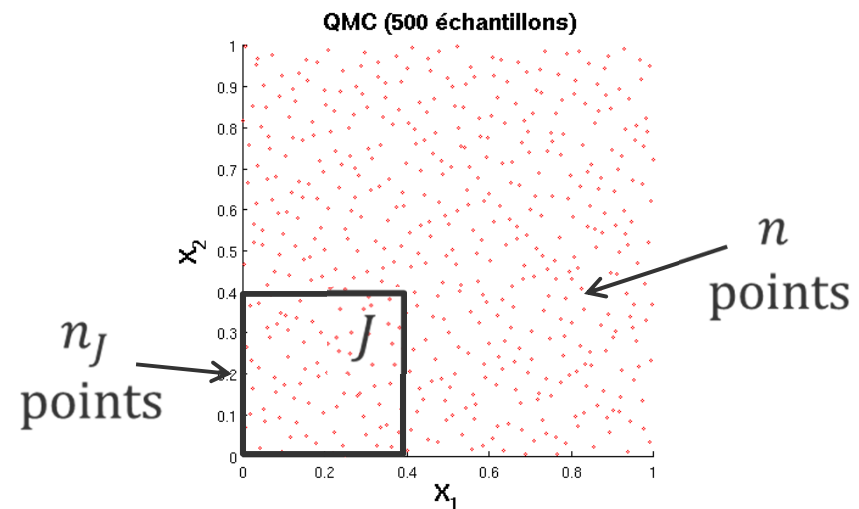
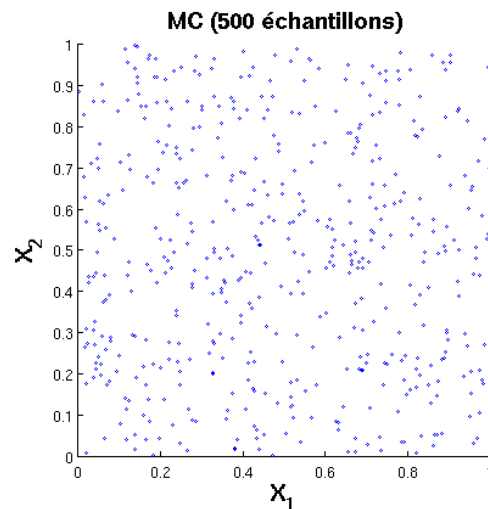
- Choix des strates  $\Omega_k$  :
  - Dans chaque strate, la variabilité de  $f$  doit être faible
- Choix du nombre de tirages  $m_k$  dans chaque strate :
  - Allocation proportionnelle  $m_k = p_k M$

$$\text{Var}[\hat{\mu}_{MC}] - \text{Var}[\hat{\mu}_{strat}] = \frac{1}{M} \sum_{k=1}^K p_k (E[f(X|Z \in \Omega_k)] - \mu)^2 > 0$$

→ Variance toujours plus faible que l'estimateur MC classique

# Quasi Monte Carlo (1)

- Améliorer l'estimation avec des tirages plus « uniformes »



- Mesure d'uniformité : la discrédance

$$D_n = \sup_J \left| \frac{n_J}{n} - Vol(J) \right|$$

Proportion de points dans  $J$       Proportion du volume occupé par  $J$

# Quasi Monte Carlo (2)

- Majoration de l'erreur commise sur l'estimation (inégalité de Koksma-Hlawka) :

$$\left| \frac{1}{n} \sum_{i=1}^n f(X_i) - E[f(X)] \right| \leq c(f) D_n(X_i)$$

Discrépance

- Suites à faible discrétance  $D_n(X_i)$  :

— Suites de Sobol, Halton, Faure, ...

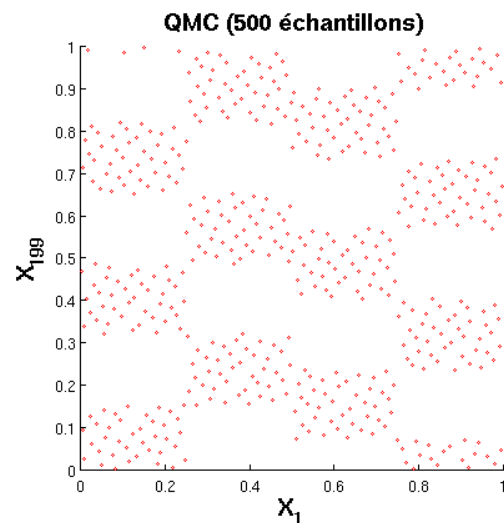
- Erreur

Dimension

$$\text{QMC} : O\left(\frac{\log(n)^d}{n}\right) < \text{MC} : O\left(\frac{1}{\sqrt{n}}\right)$$

# Quasi Monte Carlo (3)

- Suites déterministes :
  - Impossible d'appliquer le Théorème Central Limite (pour calculer des intervalles de confiance)
- Problème en grande dimension



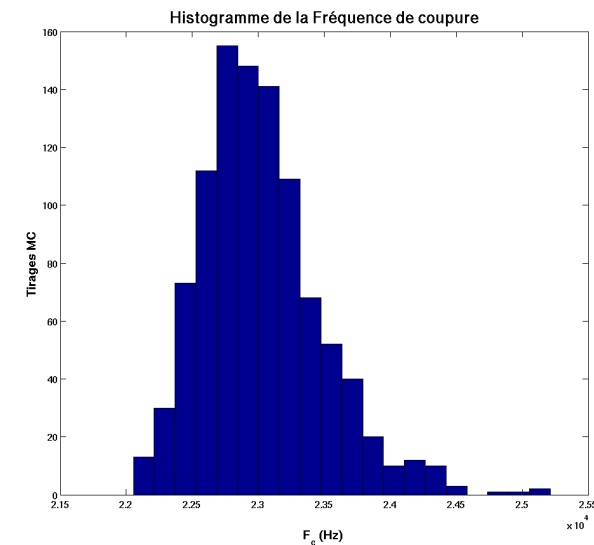
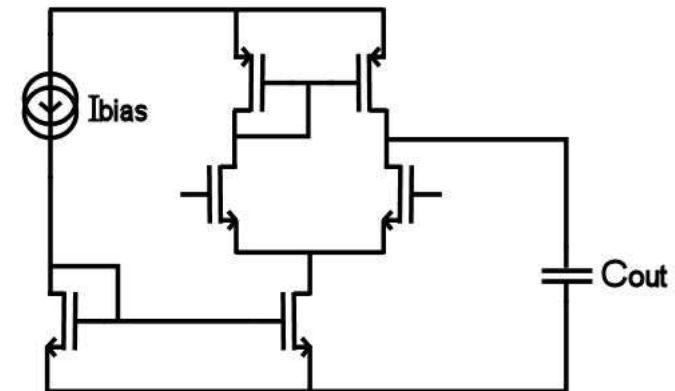
200 dimensions  
Projection  
sur  $x_1$  et  $x_{199}$

- Besoin de beaucoup de points pour assurer une bonne répartition uniforme dans l'espace



# Circuit de test – comparaison des méthodes

- Amplificateur à transconductance (CMOS 65nm)
  - 5 variations paramétriques
  - Fréquence de coupure
- Estimation de l'espérance
  - 100 tirages aléatoires
- Variance de l'estimateur
  - 50 essais
- Estimateurs
  - MC classique
  - MC préférentiel (loi normale)
  - MC stratifié (3 strates, alloc propor.)
  - QMC (suite de Halton)



# Circuit de test – comparaison des méthodes

- Résultats (100 tirages MC, 50 essais)

Méthode	Moyenne	Variance de la moyenne	Var(MC) / Var
MC classique	22 997,8 Hz	2 198,3 Hz	1
MC préférentiel	22 997,5 Hz	2 070,4 Hz	1,06
MC stratifié	23 009,8 Hz	1 417,5 Hz	1,55
QMC	22 996,0 Hz	942, 1 Hz	2,33

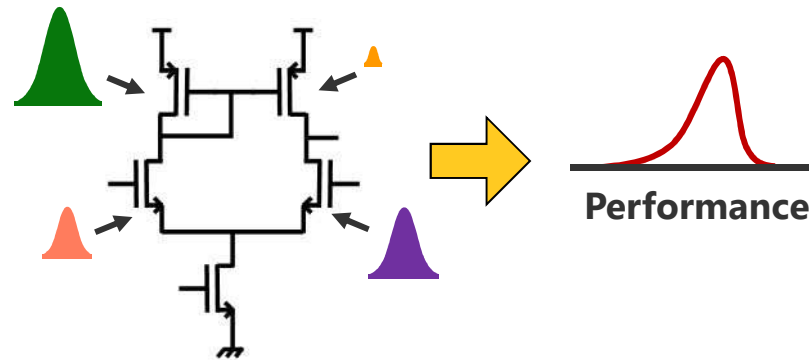
- Réduction de variance (faible)

→ Rechercher les paramètres optimaux des méthodes de réduction de variance

# MÉTHODES D'ANALYSE DE SENSIBILITÉ

# Méthodes d'analyse de la sensibilité

- Impact des variations paramétriques sur une performance du circuit ?



- Objectifs :
  - Identifier les composants du circuits dont les variations influencent le plus la performance
  - Modélisation/Simulation avec un jeu réduit de variations paramétriques (QMC)

# Analyse de sensibilité en grandes dimensions

## ■ Contexte

— ↗ du nombre de variations paramétriques ( $p = 100 \sim 1000$ )

## ■ Problématique

— Coût CPU élevé, nombre limité de simulations MC :  $p \gg M$   
→ Un nombre réduit  $p_{imp}$  de variations sont vraiment influentes

## ■ Objectif

— Modèle linéaire avec seulement les variations influentes

$$Y = \underbrace{\beta_s X_s + \dots + \beta_u X_u}_{p_{imp} \text{ termes}} = X_{imp}^T \beta_{imp}$$

## ■ Minimisation des moindres carrés

$$\min_{\beta} \sum_{k=1}^M (Y_k - \mathbf{X}_k \beta)^2 = \min_{\beta} \|Y - \mathbf{X}\beta\|^2 \rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$$

## ■ Problème : $p \leq M$

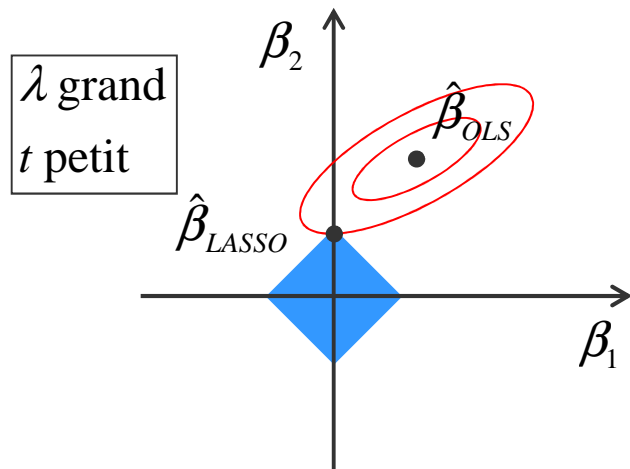
→ Solution : régularisation par pénalisation L1

# Pénalisation L1 - Régression LASSO

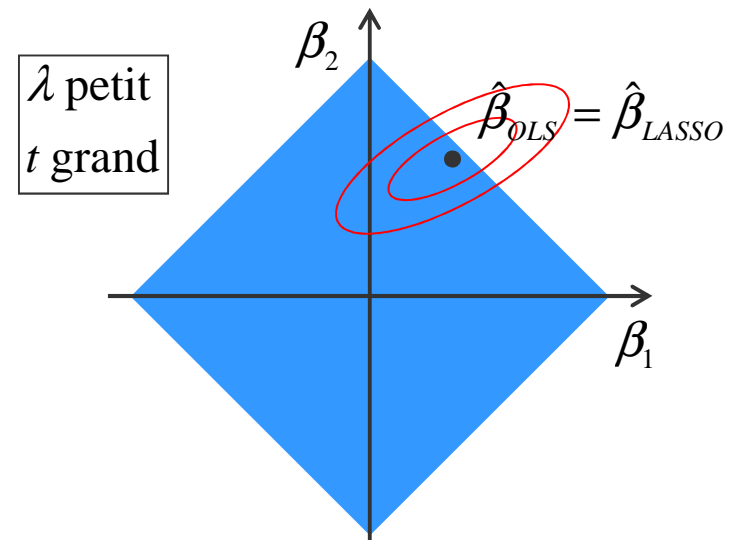
- Pénalisation L1 d'un problème de régression linéaire (LASSO)

$$\min_{\beta} \|Y - \mathbf{X}\beta\|_2^2 \Leftrightarrow \min_{\beta} \|Y - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

tel que  $\sum_{i=1}^p |\beta_i| \leq t$  ← Paramètres de régularisation  
 $t = f(\lambda)$



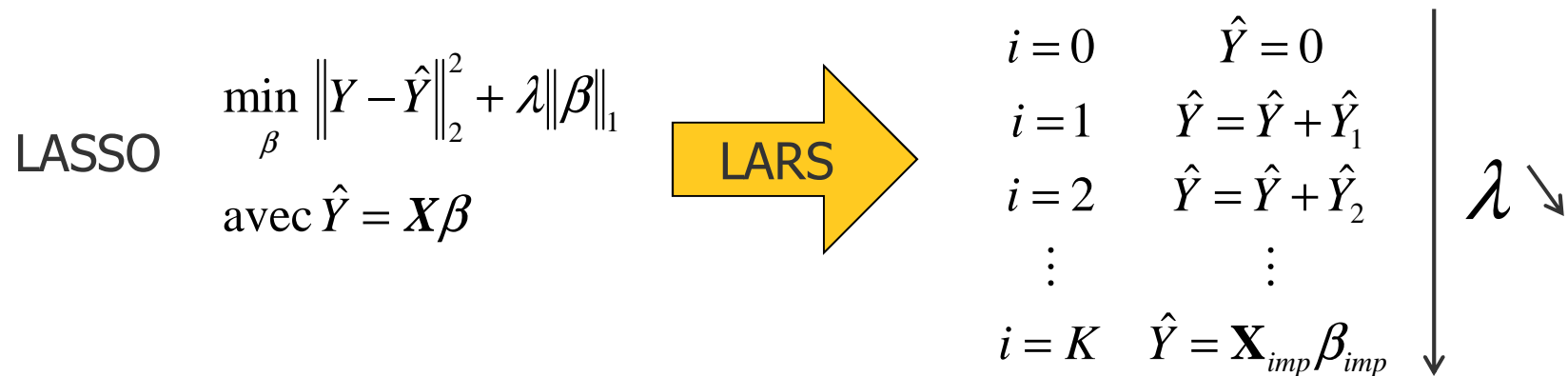
→ Seuillage de  $\beta_1$  à 0



→ Même solution que les moindres carrés

# Algorithme LARS (Least Angle regression Stagewise)

- Construction itérative du modèle linéaire
  - Ajout des variables influentes une par une



- Itération LARS

1. Trouver la variable à ajouter (la plus corrélée avec les résidus)
2. Calculer la direction de descente
3. Calculer le pas de descente

# Algorithme LARS – Conditions d'arrêt

## ■ Pas de critère d'arrêt

— Cas 1 : Nombre de variables  $\leq$  Nombre de simulations MC ( $p \leq M$ )

→ Modèle linéaire complet  $Y = \underbrace{\beta_s X_s + \dots + \beta_u X_u}_{p \text{ termes}}$  Influence décroissante  
(= moindres carrés)

— Cas 2 : Nombre de variables  $>$  Nombre de simulations MC ( $p > M$ )

→ Modèle linéaire réduit  $Y = \underbrace{\beta_s X_s + \dots + \beta_u X_u}_M$  Influence décroissante

## ■ Seuil sur les résidus

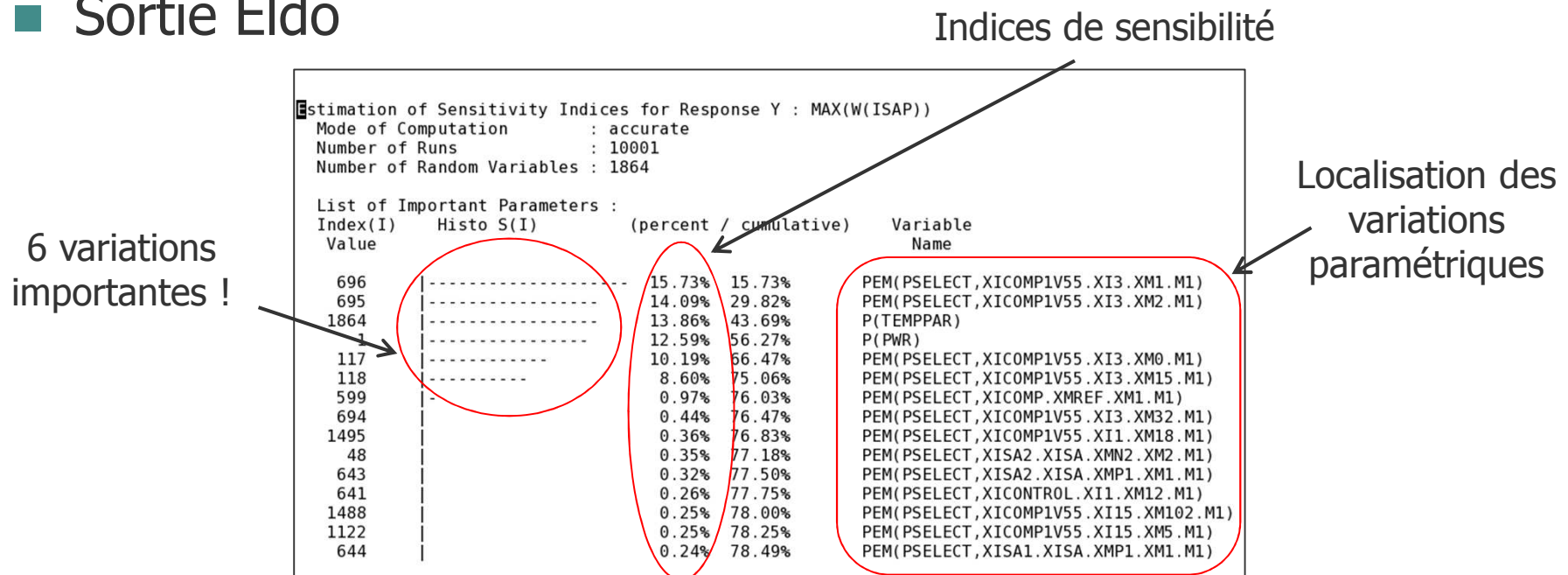
— Cas 3 :  $\left( \|Y - \hat{Y}\|_2^2 < \gamma \right)$

→ Modèle linéaire réduit  $Y = \underbrace{\beta_s X_s + \dots + \beta_u X_u}_{p_{imp} \text{ termes}}$  Influence décroissante



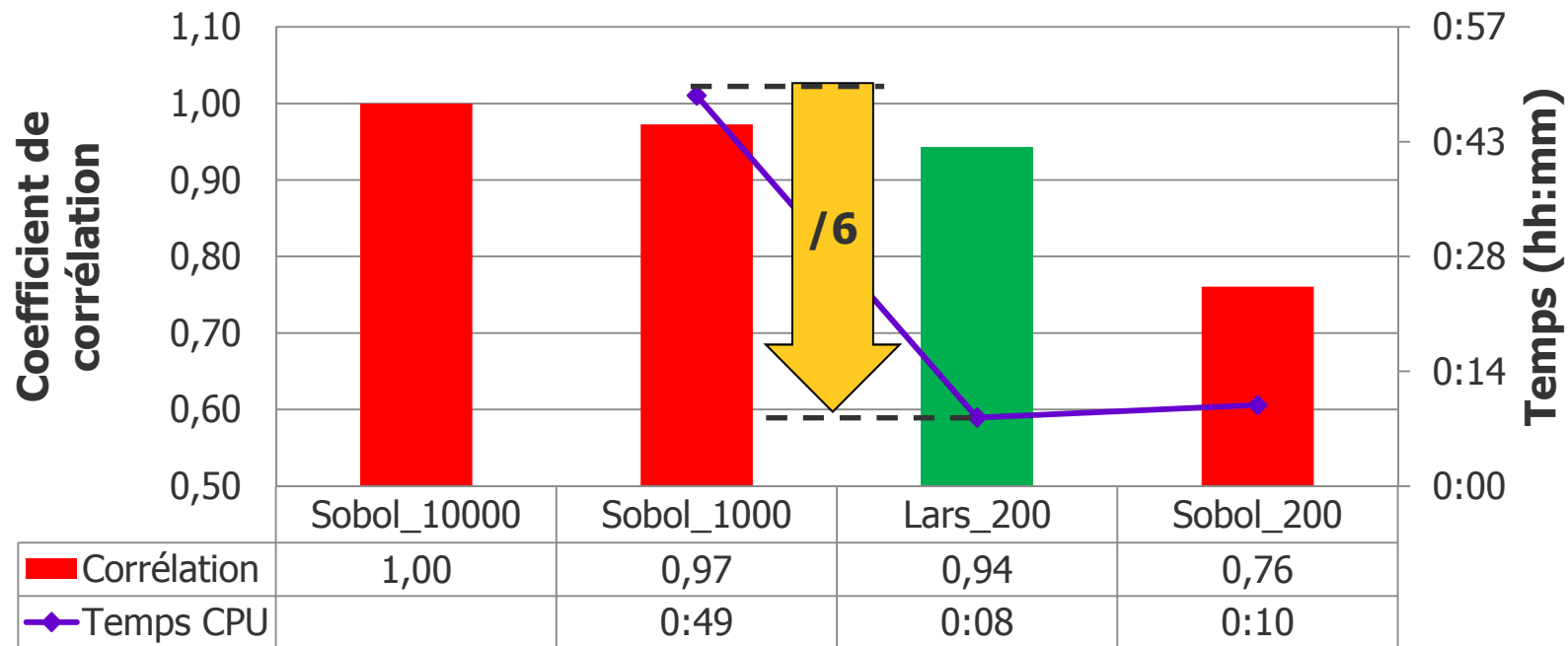
# Analyse de sensibilité de Sobol vs LARS (1)

- Circuit de test
  - 1864 variations paramétriques
- Analyse de sensibilité de référence
  - Méthode de Sobol (10000 simulations Monte Carlo, 8h09min)
- Sortie Eldo



# Analyse de sensibilité de Sobol vs LARS (2)

- Mesures pour différentes tailles d'échantillons MC
  - Coefficient de corrélation entre les indices de sensibilités
  - Temps CPU

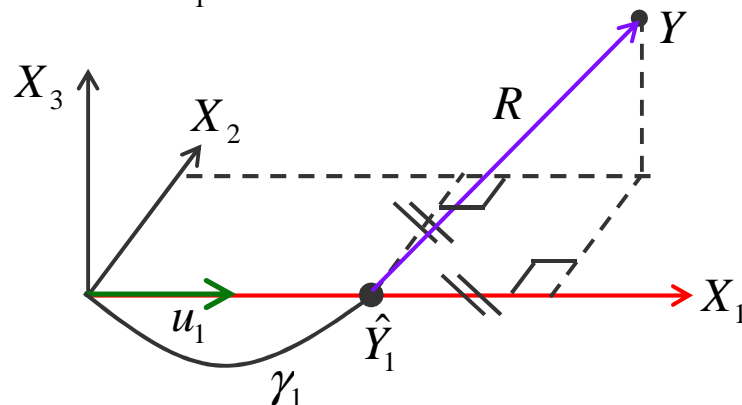
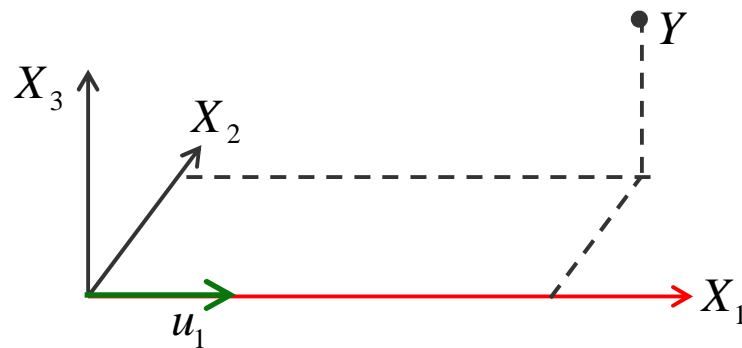
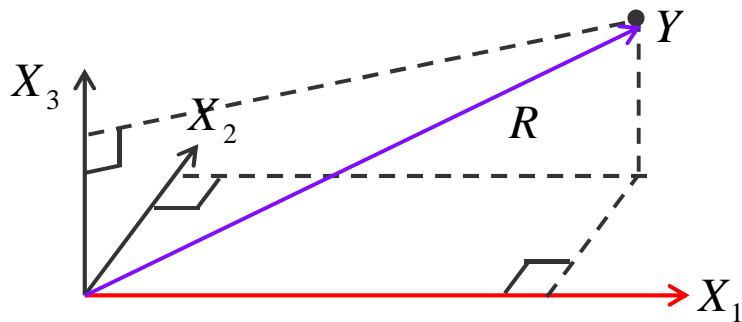


→ Lars\_200 aussi précis que Sobol\_1000 mais temps CPU divisé par 6

**Mentor  
Graphics®**

**www.mentor.com**

# Algorithme LARS – Itération 1



- Recherche de la variable la plus corrélée avec les résidus :

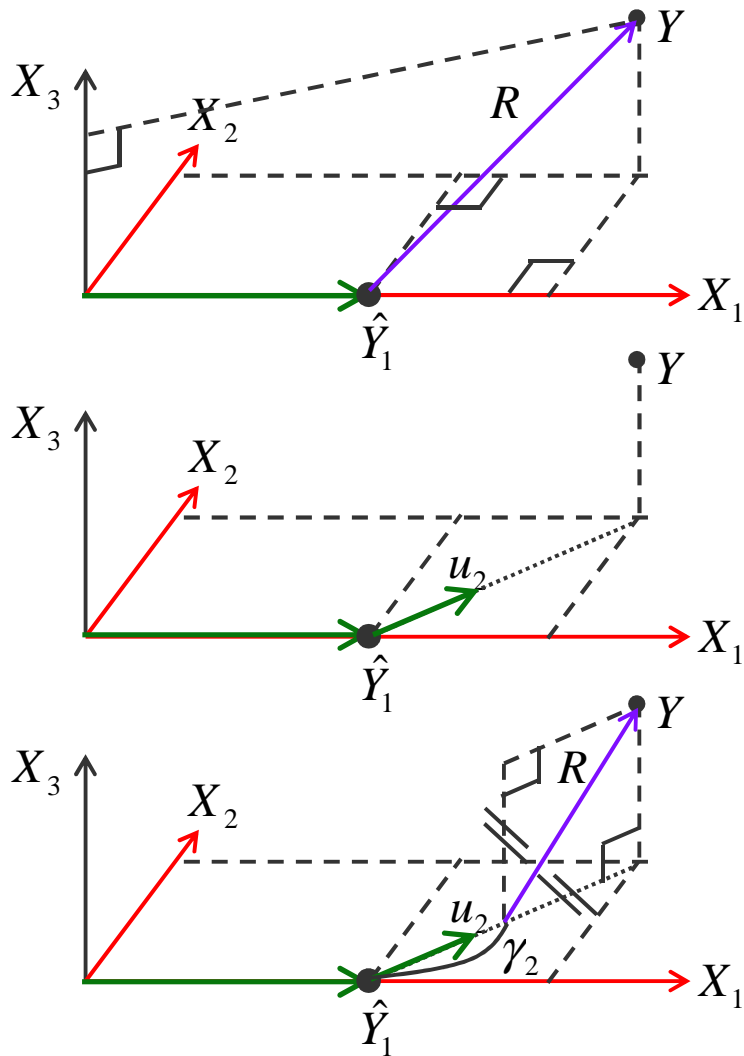
$$X^T(Y - \hat{Y}) = X^T R \rightarrow X_1$$

- Calcul de la direction  $u_1$

- Calcul du pas  $\gamma_1$ 
  - Valeur max jusqu'à ce qu'une variable inactive soit autant corrélée avec les résidus

$$\hat{Y} = \gamma_1 u_1$$

# Algorithme LARS – Itération 2



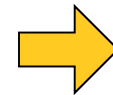
- Recherche de la variable la plus corrélée avec les résidus :

$$X^T(Y - \hat{Y}) = X^T R \rightarrow X_2$$

- Calcul de la direction  $u_2$ 
  - Vecteur équiangulaire entre les variables actives

- Calcul du pas  $\gamma_2$ 
  - Valeur max jusqu'à ce qu'une variable inactive soit autant corrélée avec les résidus

$$\hat{Y} = \gamma_1 u_1 + \gamma_2 u_2$$



Progression = compromis entre les variables les plus corrélées avec les résidus