# Comparison of Latin Hypercube and Quasi Monte Carlo Sampling Techniques

*Sergei Kucherenko[2], Andrea Saltelli[1], Daniel Albrecht[3]*

*[1]SVT - University of Bergen (UIB) and ICTA -Universitat Autonoma de Barcelona (UAB)*
*[2]Imperial College London, UK*
*[3]The European Commission, Joint Research Centre, ISPRA(VA),  ITALY*

# Outline

Monte Carlo integration methods

Latin Hypercube sampling design

Quasi Monte Carlo methods. Sobol' sequences and their properties

Comparison of sample distributions generated by different techniques

Global Sensitivity Analysis and Effective dimensions

Comparison results

# Monte Carlo integration methods

$$I[f] = \int_{H^n} f(\vec{x}) d\vec{x}$$

see as an expectation: $I[f] = E[f(\vec{x})]$

Monte Carlo : $I_N[f] = \dfrac{1}{N} \sum_{i=1}^{N} f(\vec{z}_i)$

$\{\vec{z}_i\} -$ is a sequence of random points in $H^n$

Error: $\varepsilon = \left| I[f] - I_N[f] \right|$

$$\varepsilon_N = ( E(\varepsilon^2))^{1/2} = \frac{\sigma(f)}{N^{1/2}} \rightarrow$$

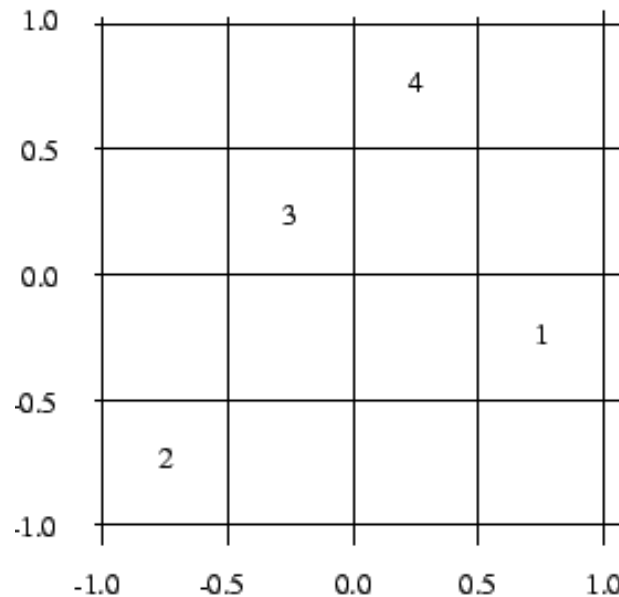Convergence does not depent on dimensionality but it is slow

Improve MC convergence by decreasing $\sigma(f)$

Use variance reduction techniques:

antithetic variables; control variates;

stratified sampling $\rightarrow$ LHS sampling

# Latin Hypercube sampling



Latin Hypercube sampling is a type of Stratified Sampling.

To sample N points in d-dimensions

Divide each dimension in N equal intervals => $N^n$ subcubes.

Take one point in each of the subcubes so that being projected to lower dimensions points do not overlap

# Latin Hypercube sampling

$\{\pi_k\}, \ k = 1, ..., n$ - independent random permutations of $\{1, ..., N\}$
each uniformly distributed over all N! possible permutations

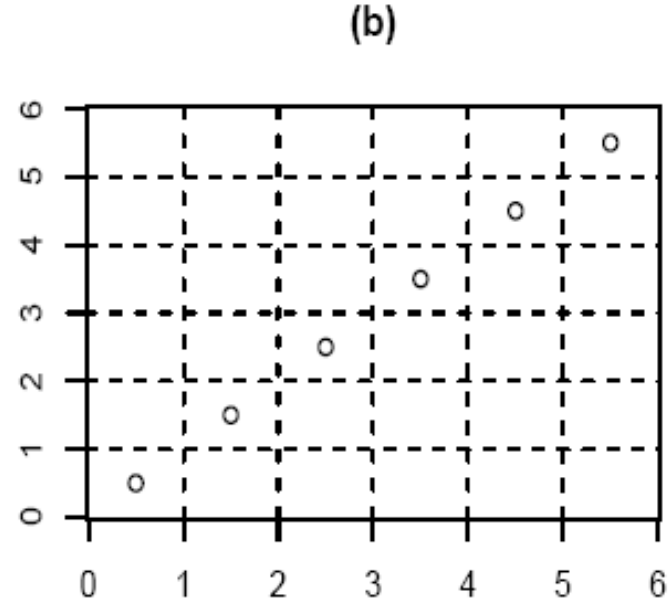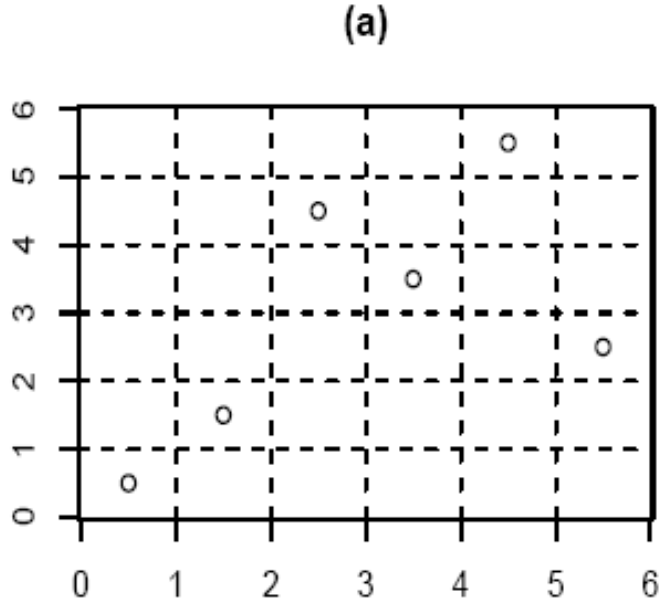LHS coordinates: $\quad x_i^k = \dfrac{\pi_k(i) - 1 + U_i^k}{N}, \ i = 1, ..., N, \ k = 1, ..., n$

$$U_i^k \approx U(0,1)$$

LHS is built by superimposing well stratified one-dimensional samples.

It cannot be expected to provide good uniformity properties in a n-dimensional unit hypercube.

# Deficiencies of LHS sampling



(a)     (b)

1) Space is badly explored (a)

2) Possible correlation between variables (b)

3) Points can not be sampled sequentially

=> Not suited for integration

# Discrepancy. Quasi Monte Carlo.

Discrepancy is a measure of deviation from uniformity:

Defintions: $Q(\vec{y}) \in H^n$, $Q(\vec{y}) = [0, y_1) \times [0, y_2) \times ... \times [0, y_n)$,

$m(Q)$ − volume of $Q$

$$D_N^* = \sup_{Q(\vec{y}) \in H^n} \left| \frac{N_{Q(\vec{y})}}{N} - m(Q) \right|$$

Random sequences: $D_N^* \rightarrow (\ln \ln N)/N^{1/2} \sim 1/N^{1/2}$

$$D_N^* \leq c(d) \frac{(\ln N)^n}{N}$$ − Low discrepancy sequences (LDS)

Convergence: $\varepsilon_{QMC} = \left| I[f] - I_N[f] \right| \leq V(f) D_N^*$,

$$\varepsilon_{QMC} = \frac{O(\ln N)^n}{N}$$

Assymptotically $\varepsilon_{QMC} \sim O(1/N) \rightarrow$ much higher than

$\varepsilon_{MC} \sim O(1/\sqrt{N})$

# QMC. Sobol' sequences

Convergence: $\varepsilon = \dfrac{O(\ln N)^n}{N}$ – for all LDS

For Sobol' LDS: $\varepsilon = \dfrac{O(\ln N)^{n-1}}{N}$, if $N = 2^k$, $k$ – integer

Sobol' LDS:

1. Best uniformity of distribution as N goes to infinity.

2. Good distribution for fairly small initial sets.

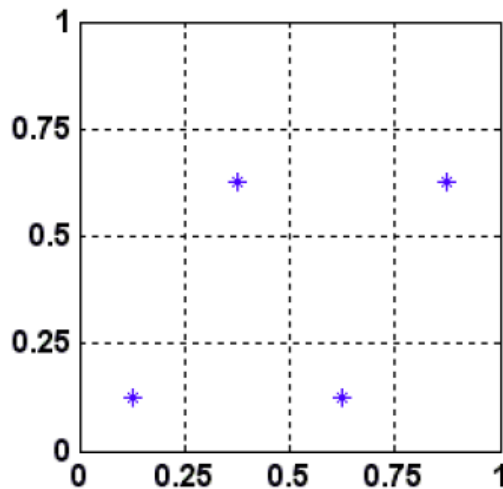3. A very fast computational algorithm.

*"Preponderance of the experimental evidence amassed to date points to Sobol' sequences as the most effective quasi-Monte Carlo method for application in financial engineering."*

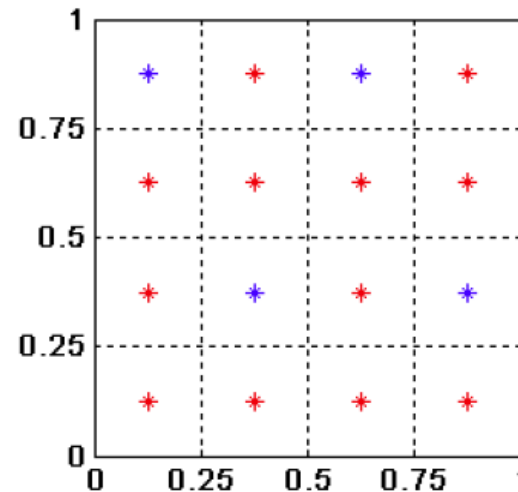Paul Glasserman, Monte Carlo Methods in Financial Engineering, Springer, 2003

A low-discrepancy sequence is said to satisfy Property A if for any binary segment (not an arbitrary subset) of the $n$-dimensional sequence of length $2^n$ there is exactly one point in each $2^n$ hyper-octant that results from subdividing the unit hypercube along each of its length extensions into half.

A low-discrepancy sequence is said to satisfy Property A' if for any binary segment (not an arbitrary subset) of the $n$-dimensional sequence of length $4^n$ there is exactly one point in each $4^n$ hyper-octant that results from subdividing the unit hypercube along each of its length extensions into four equal parts.
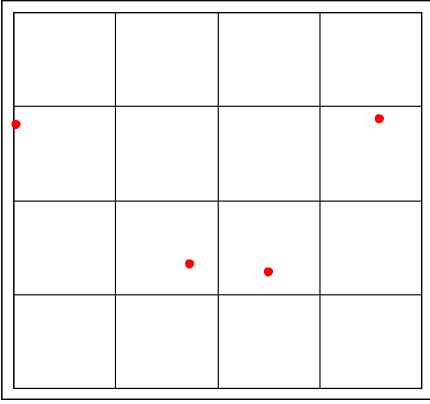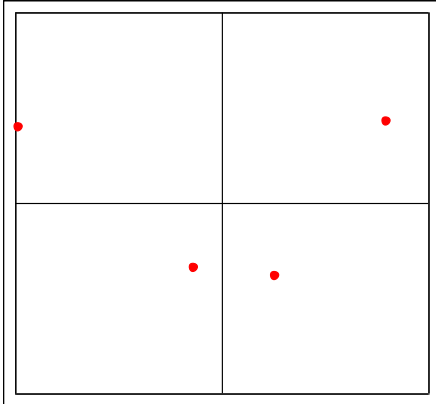


Property A                                      Property A'
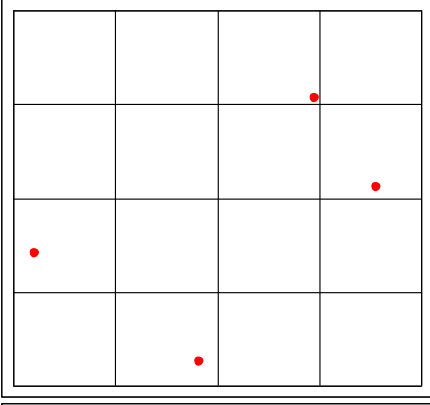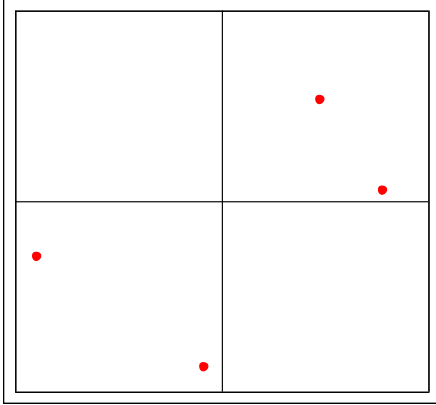
# Distributions of 4 points in two dimensions
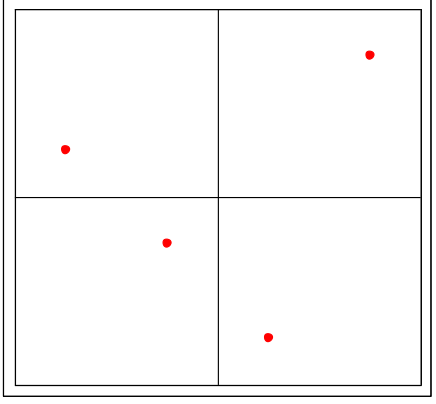
**MC ->**

**LHS ->**

**Sobol' ->**

**Property A**

**No**

**No**

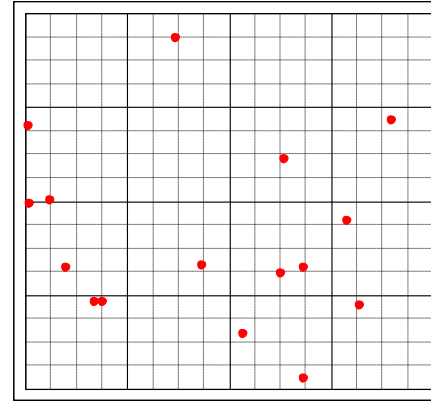**Yes**

# Distributions of 16 points in two dimensions

**MC ->**

**LHS ->**

**Sobol' ->**

**Property A'**

**No**

**No**

**Yes**

# Comparison of Discrepancy I.
# Low Dimensions



*Use standard MC and ,*

*LHS generators*

Sobol' sequence  generator:

SobolSeq:

Sobol' sequences satisfy

Properties A  and A'

www.broda.co.uk

*Result:*

*QMC in low dimensions shows much smaller discrepancy than MC and LHS*

12

# ANOVA decomposition and Sensitivity Indices

Consider a model

x is a vector of input variables

f(x) is integrable

$$Y = f(x)$$
$$x = (x_1, x_2, ..., x_k)$$
$$0 \leq x_i \leq 1$$

ANOVA decomposition:

$$Y = f(x) = f_0 + \sum_{i=1}^{k} f_i(x_i) + \sum_{i} \sum_{j>i} f_{ij}(x_i, x_j) + ... + f_{1,2,...,k}(x_1, x_2, ..., x_k),$$

$$\int_0^1 f_{i_1 \cdots i_s}(x_{i_1}, ,..., x_{i_s}) \ dx_{i_k} = 0, \ \ \forall k, \ 1 \leq k \leq s$$

Variance decomposition:

$$\sigma^2 = \sum_i \sigma_i^2 + \sum_{i,j} \sigma_{ij}^2 + ... \sigma_{1,2,...,n}^2$$

Sobol' SI:

$$1 = \sum_{i=1}^{k} S_i + \sum_{i<j} S_{ij} + \sum_{i<j<l} S_{ijl} + ... + S_{1,2,...k}$$

# Sobol' Sensitivity Indices (SI)

- **Definition:** $\boxed{S_{i_1 \ldots i_s} = \sigma^2_{i_1 \ldots i_s} / \sigma^2}$

$$\sigma^2_{i_1 \ldots i_s} = \int_0^1 f^2_{i_1 \ldots i_s}\left(x_{i_1}, \ldots, x_{is}\right) dx_{i_1}, \ldots, x_{is}$$ **- partial variances**

$$\sigma^2 = \int_0^1 \left(f(x) - f_0\right)^2 dx$$ **- total variance**

- **Sensitivity indices for subsets of variables:** $x = (y, z)$

$$\sigma^2_y = \sum_{s=1}^{m} \sum_{(i_1 \langle \ldots \langle i_s) \in \mathrm{K}} \sigma^2_{i_1, \ldots, i_s}$$

**Total variance for a subset:** $\left(\sigma^{tot}_y\right)^2 = \sigma^2 - \sigma^2_z$

**Corresponding global sensitivity indices:**

$$S_y = \sigma^2_y / \sigma^2, \qquad S^{tot}_y = \left(\sigma^{tot}_y\right)^2 / \sigma^2.$$

# Effective dimensions

Let $|u|$ be a cardinality of a set of variables $u$.

The effective dimension of $f(x)$ in the superposition sense
is the smallest integer $d_S$ such that

$$\sum_{0<|u|<d_S} S_u \geq (1-\varepsilon), \ \varepsilon << 1$$

It means that $f(x)$ is almost a sum of $d_S$-dimensional functions.

_____

The function $f(x)$ has effective dimension in the truncation sense $d_T$ if

$$\sum_{u \subseteq \{1,2,...,d_T\}} S_u \geq (1-\varepsilon), \ \varepsilon << 1$$

Important property: $d_S \leq d_T$

Example: $f(x) = \sum_{i=1}^{n} x_i \rightarrow d_S = 1, \ d_T = n$

# Classification of functions

Type A. Variables are not equally important

$$\frac{S_y^T}{n_y} >> \frac{S_z^T}{n_z} \leftrightarrow d_T << n$$

Type B,C. Variables are equally important

$$S_i \approx S_j \leftrightarrow d_T \approx n$$

Type B. Dominant low order indices

$$\sum_{i=1}^n S_i \approx 1 \leftrightarrow d_S << n$$

Type C. Dominant higher order indices

$$\sum_{i=1}^n S_i << 1 \leftrightarrow d_S \approx n$$

ANOVA: $\quad f(x) = f_0 + \sum_i f_i(x_i) + r(x)$

$r(x) -$ high order interactions terms

LHS: $\quad E(\varepsilon^2_{LHS}) = \dfrac{1}{N} \int_{H^n} [r(x)]^2 dx + O(\dfrac{1}{N})$ $\quad$ (Stein, 1987)

MC: $\quad E(\varepsilon^2_{MC}) = \dfrac{1}{N} \sum_i \int_{H^n} [f_i(x_i)]^2 dx + \dfrac{1}{N} \int_{H^n} [r(x)]^2 dx + O(\dfrac{1}{N})$

if $\quad \displaystyle\int_{H^n} [r(x)]^2 dx$ $\;$ is small $\Leftrightarrow d_S$ ( Type B functions )

$\rightarrow \quad E(\varepsilon^2_{LHS}) < E(\varepsilon^2_{MC})$

# Classification of functions

| Function type | Description | Relationship between $S_i$ and $S_i^{tot}$ | $d_T$ | $d_S$ | QMC is more efficient than MC | LHS is more efficient than MC |
|---|---|---|---|---|---|---|
| A | A few dominant variables | $S_y^{tot}/n_y >> S_z^{to}/n_z$ | $<< n$ | $<< n$ | Yes | No |
| B | No unimportant subsets; only low-order interaction terms are present | $S_i \approx S_j, \ \forall \ i, j$ <br> $S_i / S_i^{tot} \approx 1, \ \forall \ i$ | $\approx n$ | $<< n$ | Yes | Yes |
| C | No unimportant subsets; high-order interaction terms are present | $S_i \approx S_j, \ \forall \ i, j$ <br> $S_i / S_i^{tot} << 1, \ \forall \ i$ | $\approx n$ | $\approx n$ | No | No |

# How to monitor convergence of MC, LHS and QMC calculations ?

The root mean square error is defined as

$$\varepsilon = \left( \frac{1}{K} \sum_{k=1}^{K} (I_d - I_N^k)^2 \right)^{1/2}$$

$K$ is a number of independent runs

MC and LHS: all runs should be statistically independent ( use a different seed point ).

QMC: for each run a different part of the Sobol' LDS was used ( start from a different index number ).

The root mean square error is approximated by the formula

$$cN^{-\alpha}, \ 0 < \alpha < 1$$

$$\text{MC:} \quad \alpha \approx 0.5$$

$$\text{QMC:} \ \alpha \leq 1$$
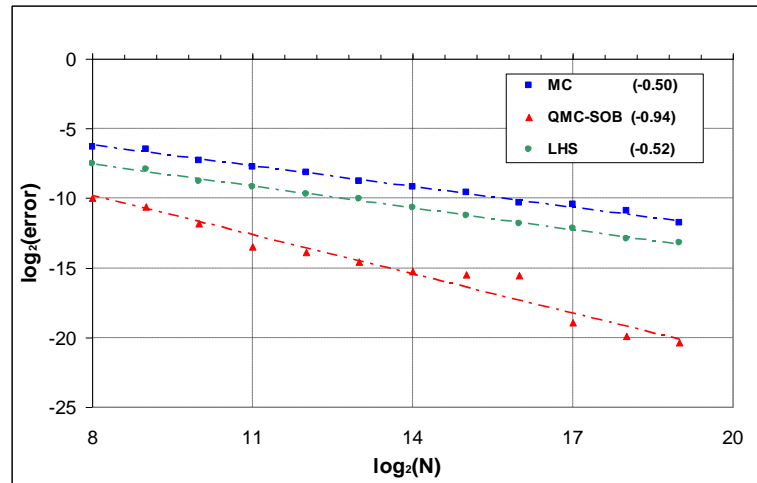
$$\text{LHS:} \ \alpha \ ?$$

# Integration error vs. N. Type A

(a) $f(x) = \sum_{j=1}^{n}(-1)^i \prod_{j=1}^{i} x_j$, $n = 360$, (b) $f(x) = \prod_{i=1}^{s} |4x_i - 2|/(1 + a_i)$, $n = 100$
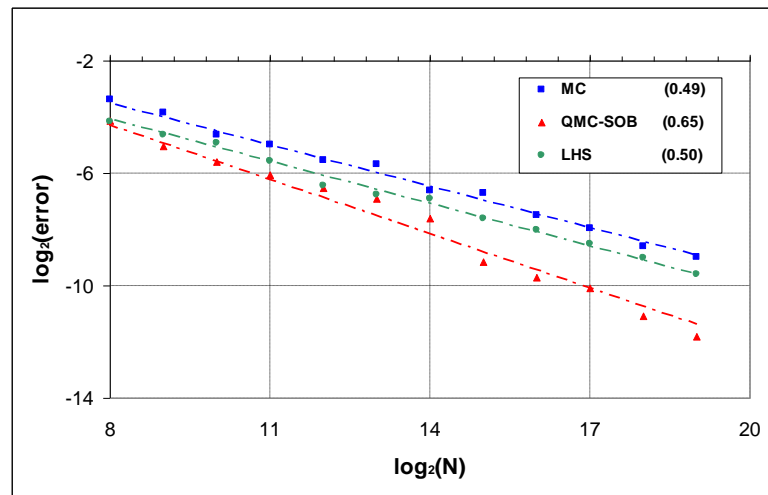
$$\varepsilon = \left( \frac{1}{K} \sum_{k=1}^{K} (I - I_N^k)^2 \right)^{1/2}$$

$$\frac{S_y^T}{n_y} >> \frac{S_z^T}{n_z} \leftrightarrow d_T << n$$

$$\varepsilon \sim N^{-\alpha}, \ 0 < \alpha < 1$$



(a)

(b)

# Integration error. Type A

$$\varepsilon = \left( \frac{1}{K} \sum_{k=1}^{K} (I - I_N^k)^2 \right)^{1/2}$$

$$\frac{S_y^T}{n_y} >> \frac{S_z^T}{n_z} \leftrightarrow d_T << n$$
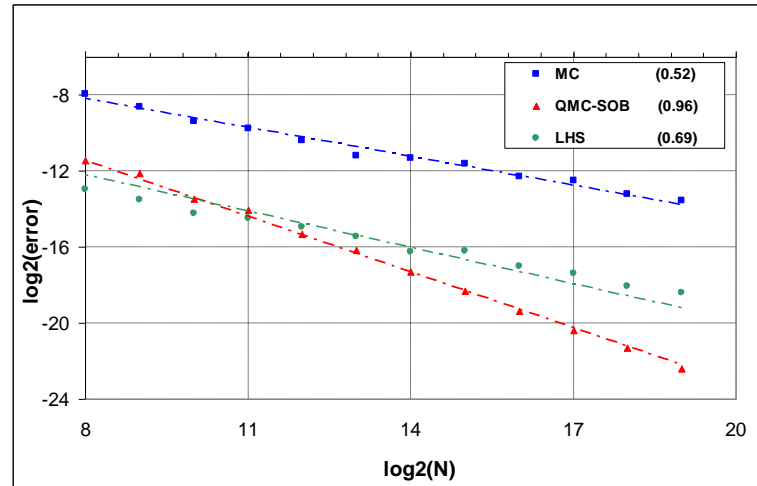
$$\varepsilon \sim N^{-\alpha}, \ 0 < \alpha < 1$$

| Index | Function | Dim $n$ | Slope MC | Slope QMC | Slope LHS |
|-------|----------|---------|----------|-----------|-----------|
| 1A | $\sum_{i=1}^{n} (-1)^i \prod_{j=1}^{i} x_j$ | 360 | 0.50 | 0.94 | 0.52 |
| 2A | $\prod_{i=1}^{n} \frac{\|4x_i - 2\| + a_i}{1 + a_i}$ <br> $a_1 = a_2 = 0$ <br> $a_3 = \ldots = a_{100} = 6.52$ | 100 | 0.49 | 0.65 | 0.50 |

# Integration error vs. N. Type B

Dominant low order indices

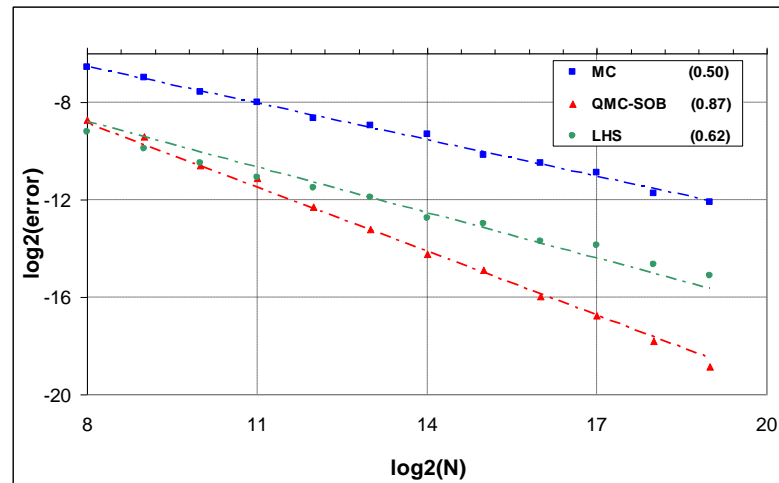$$\sum_{i=1}^{n} S_i \approx 1 \leftrightarrow d_S << n$$

(a)



$$f(x) = \prod_{i=1}^{n} \frac{n - x_i}{n - 0.5}$$

$$n = 360$$

(b)



$$f(x) = \prod_{i=1}^{n} (1 + 1/n) x_i^{1/n}$$

$$n = 360$$

22

# Integration error. Type B functions

Dominant low order indices

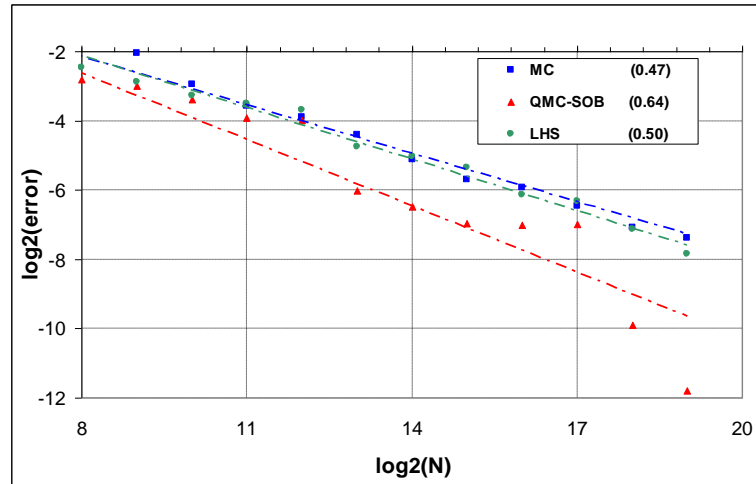$$\sum_{i=1}^{n} S_i \approx 1 \leftrightarrow d_S << n$$

| Index | Function | Dim $n$ | Slope MC | Slope QMC | Slope LHS |
|-------|----------|---------|----------|-----------|-----------|
| 1B | $\prod_{i=1}^{n} \dfrac{n - x_i}{n - 0.5}$ | 30 | 0.52 | 0.96 | 0.69 |
| 2B | $\left(1 + \dfrac{1}{n}\right)^n \prod_{i=1}^{n} \sqrt[n]{x_i}$ | 30 | 0.50 | 0.87 | 0.62 |
| 3B | $\prod_{i=1}^{n} \dfrac{\left|4 x_i - 2\right| + a_i}{1 + a_i}$  $a_i = 6.52$ | 30 | 0.51 | 0.85 | 0.55 |

Dominant higher order indices:

$$\sum_{i=1}^{n} S_i << 1 \leftrightarrow d_S \approx n$$
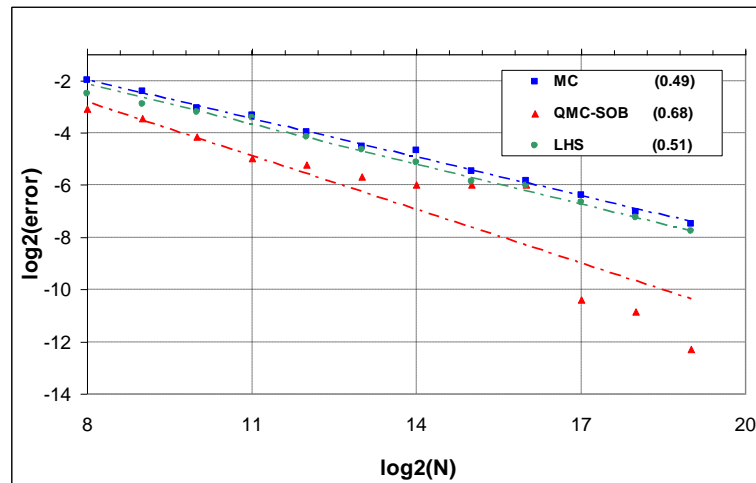
(a)



$$f(x) = \prod_{i=1}^{n} \frac{|4x_i - 2| + a_i}{1 + a_i}, a_i = 0$$

$$\rightarrow \prod_{i=1}^{n} |4x_i - 2|$$

$$n = 10$$

(b)



$$f(x) = (1/2)^{1/n} \prod_{i=1}^{n} x_i$$

$$n = 10$$

24

# Integration error for type C functions
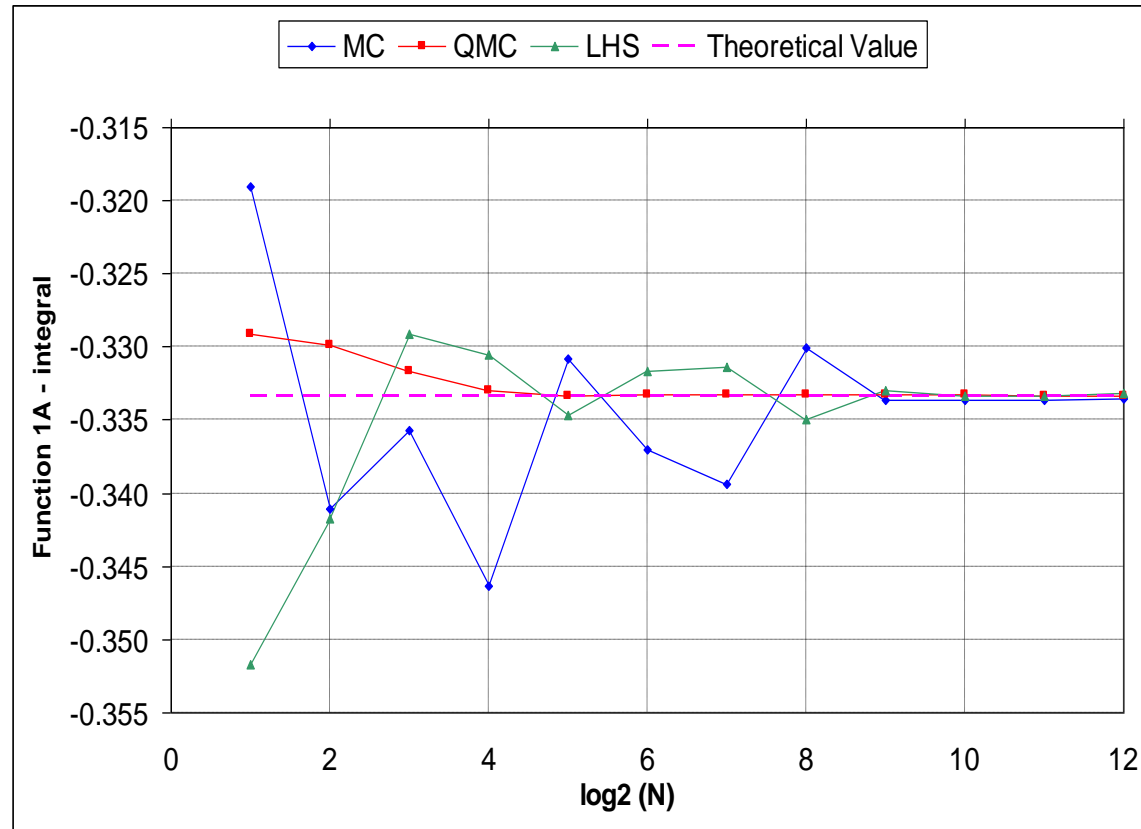
Dominant higher order indices

$$\sum_{i=1}^{n} S_i << 1 \leftrightarrow d_S \approx n$$

| Index | Function | Dim $n$ | Slope MC | Slope QMC | Slope LHS |
|-------|----------|---------|----------|-----------|-----------|
| 1C | $\prod_{i=1}^{n} \lvert 4x_i - 2 \rvert$ | 10 | 0.47 | 0.64 | 0.50 |
| 2C | $(2)^n \prod_{i=1}^{n} x_i$ | 10 | 0.49 | 0.68 | 0.51 |

# The integration error vs. N. Function 1A

$$\sum_{i=1}^{n}(-1)^{i}\prod_{j=1}^{i}x_{j},$$

$$n = 360$$



QMC: convergence is monotonic
MC and LHS: convergence curves are oscillating

QMC is 30 times faster than MC and LHS

LHS: it is not possible to incrementally add a new point while keeping the old LHS design

# Summary

Sobol' sequences possess additional uniformity properties which MC and LHS techniques do not have (Properties A and A').

Comparison of $L_2$ discrepancies shows that the QMC method has the lowest discrepancy in low dimensions ( up to 20).

QMC method outperforms MC and LHS for types A and B functions (problems with low effective dimensions)

# Summary

LHS never outperforms QMC. LHS method outperforms MC only for type B functions.

QMC remains the most efficient method among the three techniques for non-uniform distributions .

QMC should be preferred as
      better theoretical properties (A, A')
      More important variables can be associated to leftmost columns
      Sequences can be extended (automated stopping rules)
      Sequences can be replicated exactly