

Quasi-Monte Carlo

Art B. Owen
Stanford University

A tutorial reflecting personal views on what is important in QMC.

MCQMC 2016 will be at Stanford, August 14-19
mcqmc2016.stanford.edu

These are the slides I presented at MCMSki 5 in Lenzerheide Switzerland on Thursday January 8. I have added a few interstitial slides like this one after the fact. I correct some typos too.

It was my good fortune that the tutorial I gave was opposite the Imposteriors rehearsal in the other room. That gave the tutorial a near-plenary audience. While I regret not reaching the Imposteriors as well, I count myself lucky for having such good attendance.

What I present is a statistician's view of QMC, and so my emphases are different in some ways from the customary presentation of QMC.

Software Matlab has an implementation of scrambled Sobol' and Halton points. Start with the documentation for `sobolset` and `haltonset`. They support skipping initial points and leaping (taking every k 'th point for $k > 1$). I know of no reason to ever use leaping. Skipping over b^M values might be helpful to get non-overlapping QMC sets of size b^m for $m \leq M$. Skipping can also improve Halton points, but then again, randomizing without skipping should be as good or better.

Outline

- 1) What QMC is (stratification on steroids)
- 2) Why it works (discrepancy, variation and Koksma-Hlawka)
- 3) Digital constructions: van der Corput, Halton and nets
- 4) Lattices
- 5) Randomized QMC
- 6) Curse of dimension, tractability, weighted spaces
(room for Bayesian thinking here)
- 7) QMC \cap Particles
- 8) QMC \cap MCMC

But first: orientation

There are landmark papers where Monte Carlo was introduced/adapted for some specific problems. Here is a subset of examples.

- Physics [Metropolis et al. \(1953\)](#)
- Chemistry (reaction equations) [Gillespie \(1977\)](#)
- Financial valuation [Boyle \(1977\)](#)
- Resampling [Efron \(1979\)](#)
- OR (discrete event simulation) [Tocher & Owen \(1960\)](#)
- Bayes (maybe 5 landmarks in early days)
- Nonsmooth optimization [Kirkpatrick et al. \(1983\)](#)
- Computer graphics (path tracing) [Kajiya \(1988\)](#)

Quasi-Monte Carlo

QMC is used in plain quadrature.

Particle transport methods in physics [Jerome Spanier++](#)

Financial valuation, some early examples [Paskov & Traub 1990s](#)

Graphical rendering [Alex Keller++](#)

(They got an Oscar!)

Solving PDEs [Frances Kuo, Christoph Schwab++, 2015](#)

Particle methods [Chopin & Gerber \(2015\)](#)

What next

I expect that there are undiscovered landmark applications of QMC, where somebody applies/adapts/extends it for new problems, especially

- machine learning
- Bayes
- uncertainty quantification

The next slide presents a spectrum of methods to suit problems ranging from those so easy that they are considered solved to those so hard that they may never be solved. Things at the top can be embedded in your toaster. Things at the bottom may be nearly impossible or may require a team of PhDs to operate. There are no PhDs inside your toaster.

This is a riff on a comment Stephen Boyd made.

Plain QMC handles problems a bit better behaved than plain MC and delivers better results. To get at the really cutting edge problems tackled by methods at the bottom of the list may require switching 'plain QMC' to something else.

A spectrum

The harder the problem the further down this list we go.

- Everybody just knows it, e.g., 42
- We have closed form expression, e.g., $\sin(2\pi x_1) \times x_2$
- \exists exact deterministic black box, $b(\mathbf{x})$
- classic quadratures, e.g., Simpson's or Gauss rule
- plain quasi-Monte Carlo
- plain Monte Carlo
- MCMC, SMC
- ABC, aMCMC, variational MC
- \vdots
- Noone will ever know

MC and QMC

We estimate

$$\mu = \int_{[0,1]^d} f(\mathbf{x}) \, d\mathbf{x} \quad \text{by} \quad \hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i), \quad \mathbf{x}_i \in [0, 1]^d$$

In plain MC, the \mathbf{x}_i are IID $\mathbf{U}[0, 1]^d$. In QMC they're 'spread evenly'.

Non uniform

$$\mu = \int f(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x} \quad \text{and} \quad \hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\psi(\mathbf{x}_i)), \quad \mathbf{x}_i \in [0, 1]^d$$

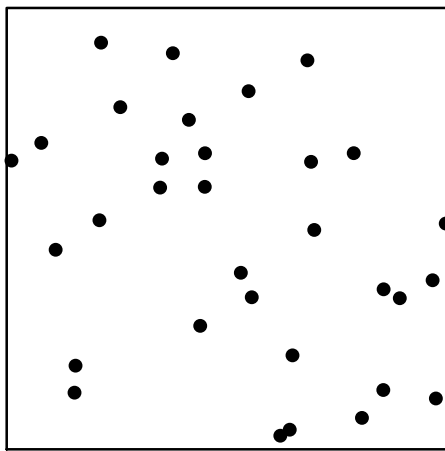
Given suitable ψ

$$\psi(\mathbf{x}) \sim p \text{ whenever } \mathbf{x} \sim \mathbf{U}[0, 1]^d$$

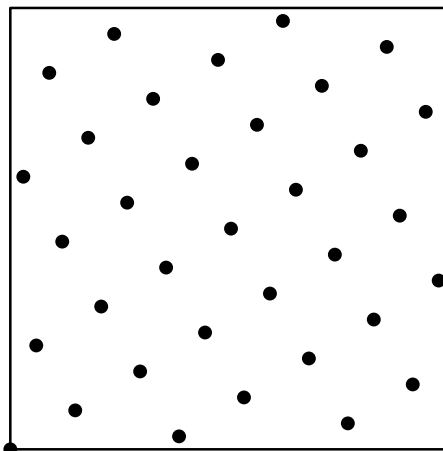
Many methods fit this framework. Acceptance-rejection is awkward.

Illustration

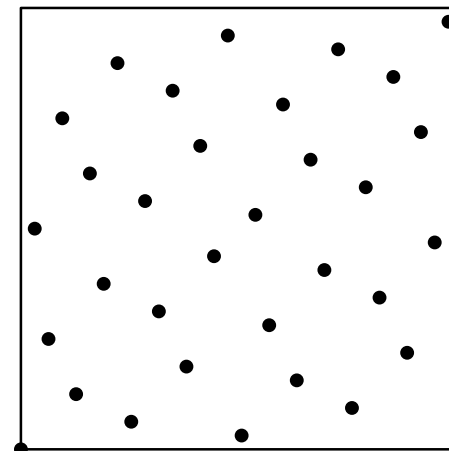
MC and two QMC methods in the unit square



Monte Carlo



Fibonacci lattice



Hammersley sequence

MC points always have clusters and gaps. What is random is where they appear.
QMC points avoid clusters and gaps to the extent that mathematics permits.

Philosophies I don't share

1) **Zaremba (1968)**: The proper justification of the normal practice of Monte Carlo integration must be based not on the randomness of the procedure, which is spurious, but on equidistribution properties of the sets of points at which the integrand values are computed.

I very much like how randomness lets you estimate your errors.

2) Also: it is possible to object to frequentist tools (LLN and CLT) being introduced into Bayesian problems.

If it works, why not use it? This complaint is now very rare.

3) Pseudo-random numbers aren't really random.

Yes, but they're very well tested. We get more trouble from floating point numbers representing reals than we do from pseudo-randomness representing randomness.

Measuring uniformity

We need a way to verify that the points \mathbf{x}_i are ‘spread out’ in $[0, 1]^d$.

The most fruitful way is to show that

$$\mathbf{U} \{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \} \doteq \mathbf{U}[0, 1]^d$$

Discrepancy

A discrepancy is a distance $\|F - \hat{F}_n\|$ between measures $F = \mathbf{U}[0, 1]^d$ and $\hat{F}_n = \mathbf{U} \{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \}$.

There are many discrepancies.

Technicality

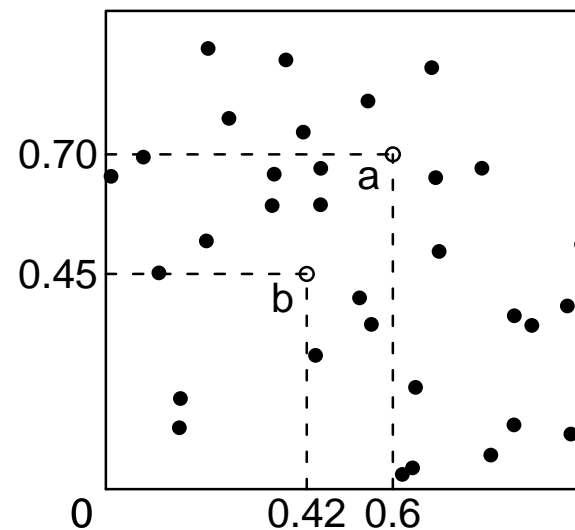
More properly: $\text{Distn}(\mathbf{x}_i) \doteq \mathbf{U}[0, 1]^d$, for $i \sim \mathbf{U}\{1, 2, \dots, n\}$

There could be ties.

Local discrepancy

Did the box $[0, \mathbf{a})$ get its fair share of points?

Local discrepancy at \mathbf{a} , \mathbf{b}



$$\delta(\mathbf{a}) = \hat{F}_n([0, \mathbf{a})) - F([0, \mathbf{a})) = \frac{13}{32} - 0.6 \times 0.7 = -0.01375$$

Star discrepancy

$$D_n^* = D_n^*(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sup_{\mathbf{a} \in [0, 1)^d} |\delta(\mathbf{a})|$$

For $d = 1$ this is Kolmogorov-Smirnov.

More discrepancies

$$D_n^* = \sup_{\mathbf{a} \in [0,1)^d} |\hat{F}_n([0, \mathbf{a})) - F([0, \mathbf{a}))|$$

$$D_n = \sup_{\mathbf{a}, \mathbf{b} \in [0,1)^d} |\hat{F}_n([\mathbf{a}, \mathbf{b})) - F([\mathbf{a}, \mathbf{b}))|$$

$$D_n^* \leq D_n \leq 2^d D_n^*$$

L^p discrepancies

$$D_n^{*p} = \left(\int_{[0,1)^d} |\delta(\mathbf{a})|^p d\mathbf{a} \right)^{1/p}$$

Also

Wrap-around discrepancies [Hickernell](#)

Discrepancies over (triangles, rotated rectangles, balls \dots convex sets \dots).

[Beck](#), [Chen](#), [Schmidt](#)

Best results are **only** for axis-aligned hyper-rectangles.

Koksma's inequality

For $d = 1$ $|\hat{\mu} - \mu| \leq D_n^*(x_1, \dots, x_n) \times \int_0^1 |f'(x)| dx$

NB: $\int_0^1 |f'(x)| dx$ is the total variation of f .

Koksma-Hlawka theorem

$$\left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) - \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} \right| \leq D_n^* \times V_{\text{HK}}(f)$$

V_{HK} is the **total variation** in the sense of [Hardy \(1905\)](#) and [Krause \(1903\)](#)

Puzzler

Is this a 100% confidence interval?

For $d = 1$, Koksma

WLOG $0 \equiv x_0 \leq x_1 \leq x_2 \leq \dots \leq x_n \leq x_{n+1} \equiv 1$

Integration and summation by parts

$$\begin{aligned} \int_0^1 f(x) \, dx &= f(1) - \int_0^1 x f'(x) \, dx \\ \frac{1}{n} \sum_{i=1}^n f(x_i) &= f(1) - \frac{1}{n} \sum_{i=0}^n i (f(x_{i+1}) - f(x_i)) \\ &= f(1) - \frac{1}{n} \sum_{i=0}^n i \int_{x_i}^{x_{i+1}} f'(x) \, dx \end{aligned}$$

After simplification, for continuous f'

$$\hat{\mu} - \mu = \frac{1}{n} \sum_{i=1}^n f(x_i) - \int_0^1 f(x) \, dx = - \int_0^1 \delta(x) f'(x) \, dx$$

Upshot

$$|\hat{\mu} - \mu| \leq \|\delta\|_p \times \|f'\|_q \quad 1/p + 1/q = 1$$

Rates of convergence

It is possible to get $D_n^* = \frac{\log(n)^{d-1}}{n}$.

Then

$$|\hat{\mu} - \mu| = o(n^{-1+\epsilon}) \quad \text{vs} \quad O_p(n^{-1/2}) \text{ for MC}$$

What about those logs?

$\log(n)^{d-1}/n$ can be very large if d is large.

The s dimensional coordinate projections (i.e., margins) of $\mathbf{x}_1, \dots, \mathbf{x}_n$ typically have discrepancy $O(\log(n)^{s-1}/n)$

Log factors are not material for functions of low effective dimension (later).

Roth (1954)

$$D_n^* \text{ of } o\left(\frac{\log(n)^{(d-1)/2}}{n}\right) \text{ is unattainable}$$

Randomization

Can also control the log factors. (later)

Tight and loose bounds

They are not mutually exclusive.

Koksma-Hlawka is tight

$$|\hat{\mu} - \mu| \leq (1 - \epsilon) D_n^*(\mathbf{x}_1, \dots, \mathbf{x}_n) \times V_{\text{HK}}(f) \quad \text{fails for some } f$$

KH holds as an equality for a worst case function like $f' \doteq \pm\delta$.

Koksma-Hlawka is also very loose

It can greatly over-estimate actual error. Usually δ and f' are dissimilar.

$$\hat{\mu} - \mu = -\langle \delta, f' \rangle$$

Just like Chebychev's inequality

It is also tight and very loose. E.g., $\Pr(|\mathcal{N}(0, 1)| \geq 10) \leq 0.01$ is loose.

Yes: $1.5 \times 10^{-23} \leq 10^{-2}$

Variation

Multidimensional variation can be a bit more subtle than one dimensional.

$$f(x_1, x_2) = \begin{cases} 1, & x_1 + x_2 \leq 1/2 \\ 0, & \text{else} \end{cases}$$

$$V_{\text{HK}}(f) = \infty \quad \text{on } [0, 1]^2$$

$$V_{\text{HK}}(f_\epsilon) < \infty, \quad \text{for } \|f - f_\epsilon\|_1 < \epsilon$$

Vitali variation

It is almost $\int_{[0,1]^d} \left| \frac{\partial^d}{\partial x_1, \dots, \partial x_d} f(\mathbf{x}) \right| d\mathbf{x}$

Vanishes if f does not depend on x_j for some $j \in 1, \dots, d$.

Hardy-Krause variation

Sum of Vitali variations on ‘upper faces’ of $[0, 1]^d$ with 0 to $d - 1$ components fixed at 1.

QMC-friendly discontinuities

Axis parallel. **Xiaoqun Wang++** $V_{\text{HK}} < \infty$.

Effective dimension

Functional ANOVA decomposition [Hoeffding \(1948\)](#), [Sobol' \(1969\)](#)

$$f(\mathbf{x}) = \mu + f_1(x_1) + f_2(x_2) + \cdots + f_d(x_d) + f_{1,2}(x_1, x_2) + \cdots$$

Each of these sub- f 's integrates to zero.

Then

$$\hat{\mu} - \mu = \frac{1}{n} \sum_{i=1}^n f_1(x_{i1}) + \cdots + \frac{1}{n} \sum_{i=1}^n f_d(x_{id}) + \frac{1}{n} \sum_{i=1}^n f_{1,2}(x_{i1}, x_{i2}) + \cdots$$

If f is dominated by its main effects and low order interactions ([Caflisch, Morokoff & O \(1997\)](#)) then the error is like a lower dimensional QMC error.

With some notational license

$$\begin{aligned} |\hat{\mu} - \mu| &\leq \sum_{u \subseteq \{1, \dots, d\}} D_n^*(\mathbf{x}_{1u}, \dots, \mathbf{x}_{nu}) V_{\text{HK}}(f_u) \\ &= \sum_{u \subseteq \{1, \dots, d\}} O\left(\frac{\log(n)^{|u|-1}}{n}\right) V_{\text{HK}}(f_u) \end{aligned}$$

For $d = 1$

Points $x_i = (i - 1/2)/n$ minimize both D_n and D_n^* .

But **where** do we put the $n + 1$ 'st point?

Extensible sequences

Take first n points of $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n, \mathbf{x}_{n+1}, \mathbf{x}_{n+2}, \dots$

Then we can get $D_n^* = O((\log n)^d/n)$.

No known extensible constructions get $O((\log n)^{d-1}/n)$.

van der Corput

i				$\phi_2(i)$
1	1	0.1	1/2	0.5
2	10	0.01	1/4	0.25
3	11	0.11	3/4	0.75
4	100	0.001	1/8	0.125
5	101	0.101	5/8	0.625
6	110	0.011	3/8	0.375
7	111	0.111	7/8	0.875
8	1000	0.0001	1/16	0.0625
9	1001	0.1001	9/16	0.5625

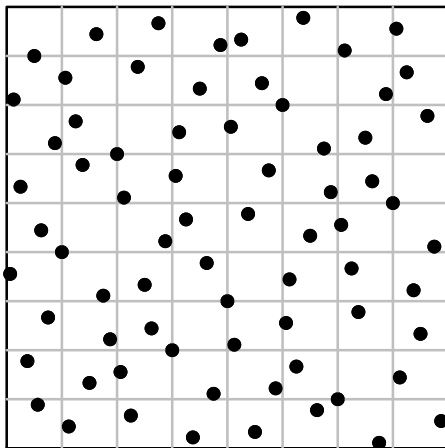
Take $x_i = \phi_2(i)$. Extensible with $D_n^* = O(\log(n)/n)$.

Commonly $x_i = \phi_2(i - 1)$ starts at $x_1 = 0$.

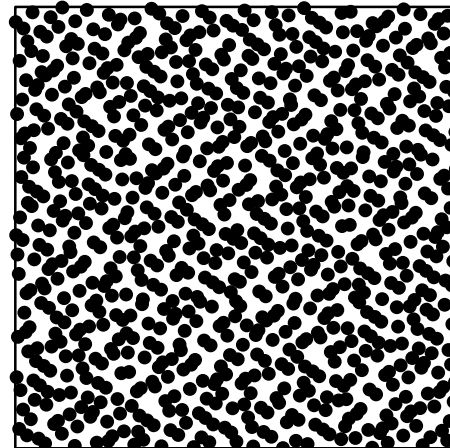
Halton sequences

The van der Corput trick works for any base. Use bases 2, 3, 5, 7, ...

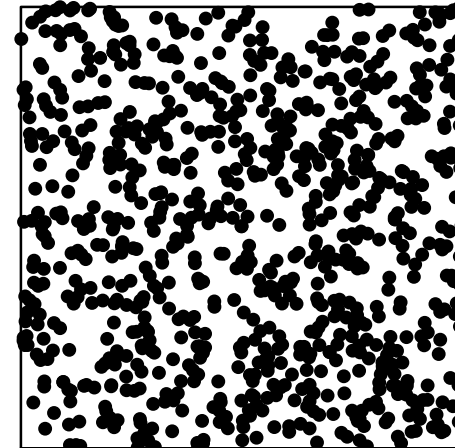
Halton sequence in the unit square



72 Halton points



864 Halton points



864 random points

Via base b digital expansions

$$i = \sum_{k=0}^{K} b^k a_{ik} \quad \rightarrow \quad \phi_b(i) \equiv \sum_{k=0}^{K} b^{-1-k} a_{ik}$$

$$\mathbf{x}_i = (\phi_2(i), \phi_3(i), \dots, \phi_p(i))$$

The Halton sequence is easy to implement. It is a good tool for testing whether QMC will help on your problem. It is also well suited to homework problems.

If you need large prime numbers, then scrambling the digits in the Halton sequence is recommended. Simply using a random permutation of them, as Giray Okten advocates, will bring an improvement.

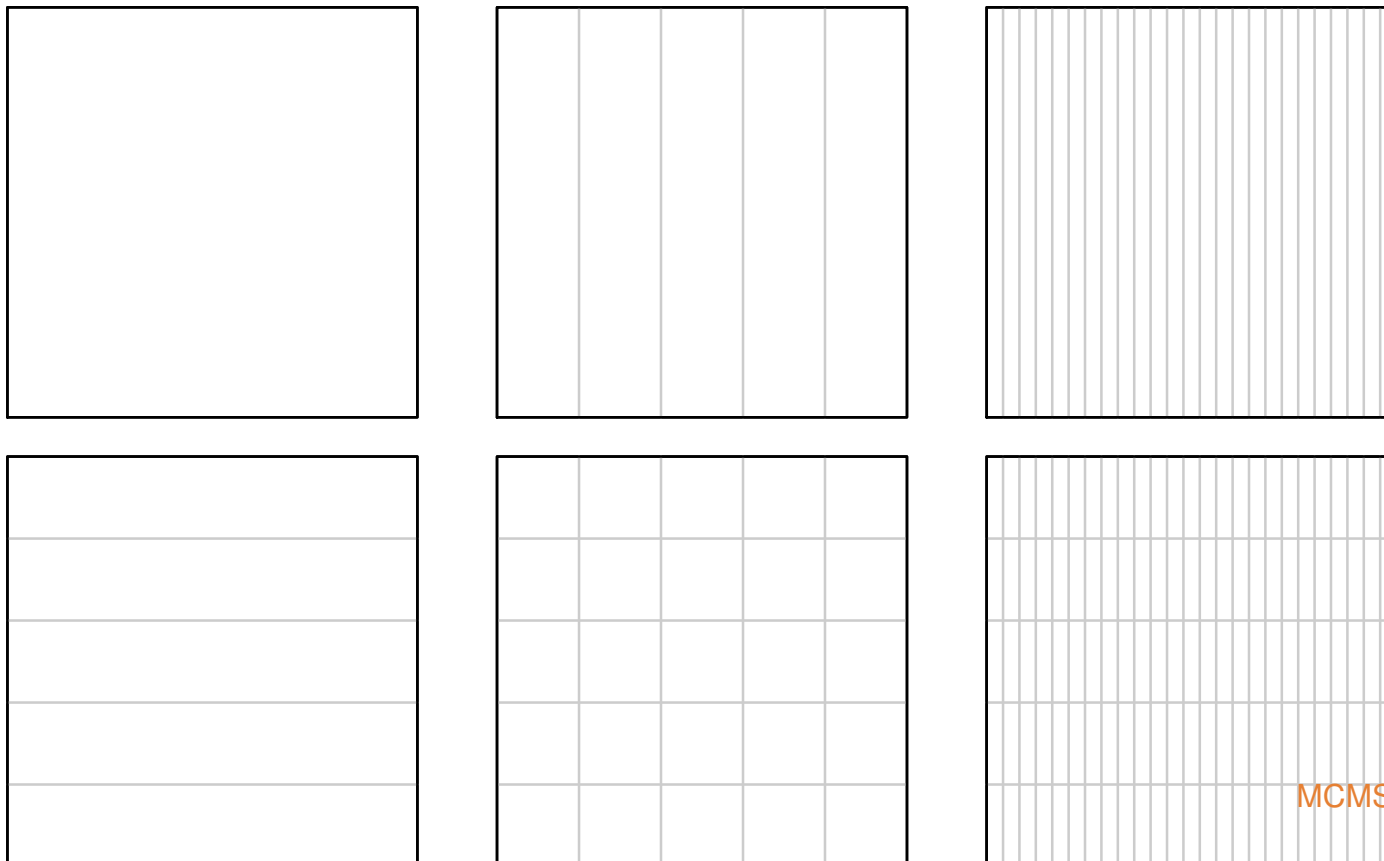
Digital nets

Halton sequences are balanced if n is a multiple of 2^a and 3^b and $5^c \dots$

Digital nets use just one base $b \implies$ balance all margins equally.

Elementary intervals

Some elementary intervals in base 5



Digital nets

$$E = \prod_{j=1}^s \left[\frac{a_j}{b^{k_j}} \frac{a_j + 1}{b^{k_j}} \right), \quad 0 \leq a_j < b^{k_j}$$

$(0, m, s)$ -net

$n = b^m$ points in $[0, 1)^s$. If $\text{vol}(E) = 1/n$ then E has one of the n points.

e.g. Faure (1982) points, prime base $b \geq s$

(t, m, s) -net

If E deserves b^t points it gets b^t points. Integer $t \geq 0$.

e.g. Sobol' (1967) points base 2

Smaller t is better (but a construction might not exist).

minT project

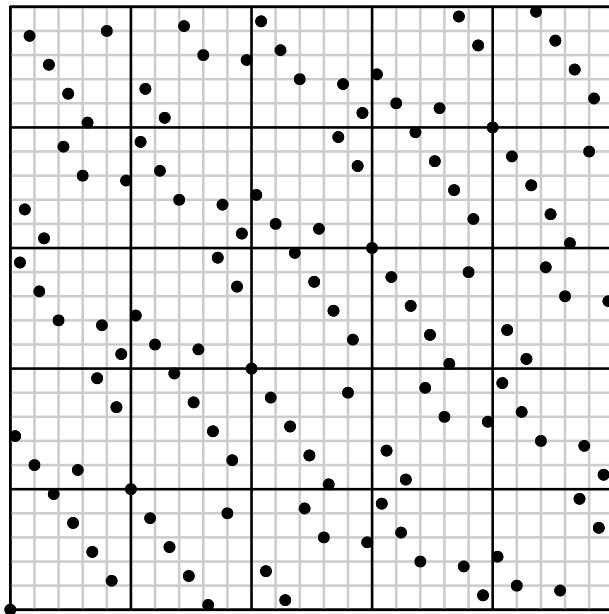
Schürer & Schmid give bounds on t given b, m and s

Monographs

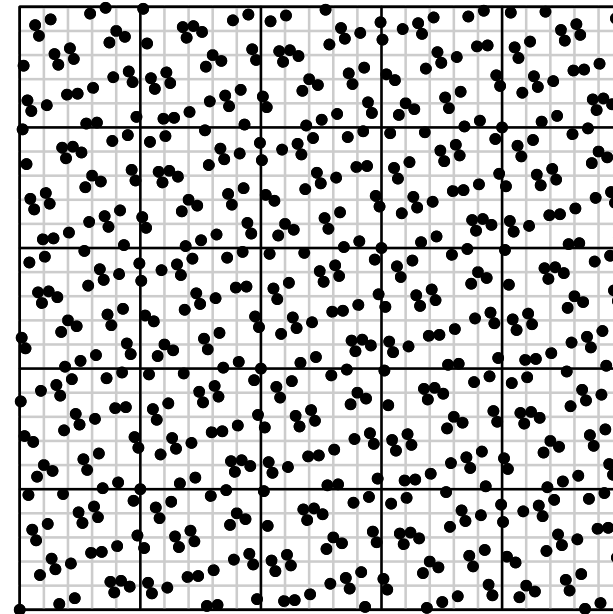
Niederreiter (1992) Dick & Pillichshammer (2010)

Example nets

Two digital nets in base 5



A (0,3,2) net



A (0,4,2) net

The $(0, 4, 2)$ -net is a bivariate margin of a $(0, 4, 5)$ -net.

The parent net has $5^4 = 625$ points in $[0, 1)^5$.

It balances 43,750 elementary intervals.

Think of 43,750 control variates for 625 obs.

We should remove that diagonal striping artifact (later).

Extensible nets

Nets can be extended to larger sample sizes.

It raises D_n^* from $O((\log n)^{s-1}/n)$ to $O((\log n)^s/n)$

(t, s) -sequence in base b

Infinite sequence of (t, m, s) -nets.

$$\underbrace{\mathbf{x}_1, \dots, \mathbf{x}_{b^m}} \quad \underbrace{\mathbf{x}_{b^m+1}, \dots, \mathbf{x}_{2b^m}} \quad \cdots \quad \underbrace{\mathbf{x}_{kb^m+1}, \dots, \mathbf{x}_{(k+1)b^m}} \quad \cdots$$

Simultaneously for all $t \geq m$.

$$\underbrace{\underbrace{(t, m, s)\text{-net}}_{1\text{st}} \quad \underbrace{(t, m, s)\text{-net}}_{2\text{nd}} \quad \cdots \quad \underbrace{(t, m, s)\text{-net}}_{b\text{'th}}}_{(t, m+1, s)\text{-net}} \quad \cdots$$

Examples

Sobol' $b = 2$ Faure $t = 0$ Niederreiter & Xing $b = 2$ (mostly)

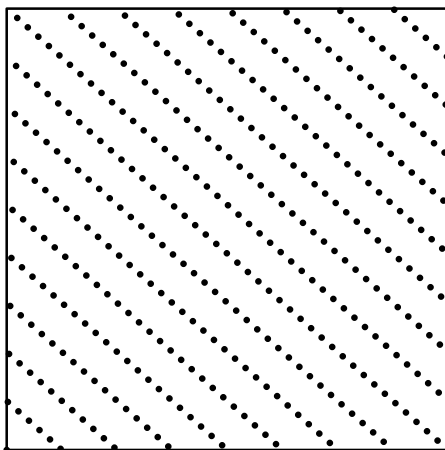
Lattices

$$\mathbf{x}_i = \left(\frac{i}{n}, \frac{Z_2 i}{n}, \frac{Z_3 i}{n}, \dots, \frac{Z_d i}{n} \right) \pmod{1} \quad Z_j \in \mathbb{N}$$

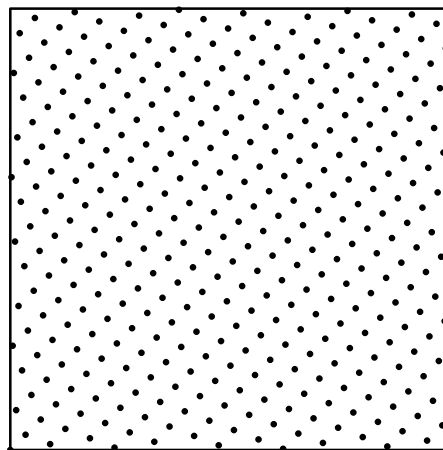
$$\mathbf{Z} = (1, Z_2, Z_3, \dots, Z_d)$$

- the other main family of QMC points
- an extensive literature, e.g., Sloan & Joe also Kuo, Nuyens, Dick, Cools, Hickernell, ...
- less benefit from randomization

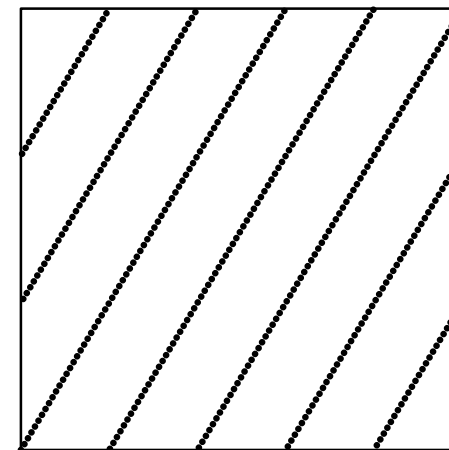
Some lattice rules for $n=377$



$z = (1, 41)$



$z = (1, 233)$



$z = (1, 253)$

Lattices ctd

The choice of $\mathbf{Z} = (1, Z_2, Z_3, \dots, Z_d) \in \mathbb{N}^d$ is critical.

We **get to** choose \mathbf{Z} . \implies We can tune/optimize.

We **have to** choose \mathbf{Z} . \implies We must search.

Korobov (1959) lattices

$$\mathbf{Z} = (1, Z, Z^2, \dots, Z^{d-1}), \quad Z \in \mathbb{N}$$

Reduced search space. Minor performance penalty.

Extensibility

Lattices are not extensible. But ‘shifted lattices’ can be extended.

Hickernell, Hong, L’Ecuyer & Lemieux (2000)

Lattice strategy

Sloan & Joe (1994) approach to lattices.

- 1) Suppose that f is periodic.
- 2) Upper bound $|\hat{\mu} - \mu|$ for smooth periodic functions.
- 3) Find lattice with small upper bound.
- 4) “Make f periodic.”

The last step involves replacing f by periodic \tilde{f} with $\int \tilde{f}(\mathbf{x}) \, d\mathbf{x} = \int f(\mathbf{x}) \, d\mathbf{x}$.

I skipped over some of these slides, to just show and explain the dual lattices a few slides hence.

This presentation of lattices is based on the book by Sloan and Joe. There has been considerable progress since then and I'm hopeful that somebody will put it together into a new monograph.

Smooth periodic functions

$$f(\mathbf{x} + \mathbf{z}) = f(\mathbf{x}), \quad \forall \mathbf{z} \in \mathbb{Z}^d$$

Fourier

$$f(\mathbf{x}) = \sum_{\mathbf{h} \in \mathbb{Z}^d} \hat{f}(\mathbf{h}) e^{2\pi\sqrt{-1}\mathbf{h}^\top \mathbf{x}} \quad (\text{in mean square})$$

$$\hat{f}(\mathbf{h}) = \int_{[0,1]^d} f(\mathbf{x}) e^{-2\pi\sqrt{-1}\mathbf{h}^\top \mathbf{x}}$$

Substitute

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) = \sum_{\mathbf{h} \in \mathbb{Z}^d} \hat{f}(\mathbf{h}) \left(\frac{1}{n} \sum_{i=1}^n e^{2\pi\mathbf{h}^\top \mathbf{x}_i} \right)$$

$$\mu = \int f(\mathbf{x}) \, d\mathbf{x} = \sum_{\mathbf{h} \in \mathbb{Z}^d} \hat{f}(\mathbf{h}) \int_{[0,1]^d} e^{2\pi\mathbf{h}^\top \mathbf{x}} \, d\mathbf{x} = \hat{f}(0)$$

Averaging sinusoids

$$\mathbf{x}_i = \frac{i\mathbf{Z}}{n} \bmod 1$$

$$\frac{1}{n} \sum_{i=1}^n e^{2\pi i \mathbf{h}^\top \mathbf{x}_i} = \begin{cases} 1 & \mathbf{h}^\top \mathbf{Z} = 0 \pmod{n} \\ 0 & \text{else.} \end{cases}$$

$$\hat{\mu} - \mu = \sum_{\substack{\mathbf{h} \in \mathbb{Z}^d \\ \mathbf{h}^\top \mathbf{Z} = 0 \pmod{n}}} \hat{f}(\mathbf{h}) - \hat{f}(\mathbf{0})$$

Dual lattices

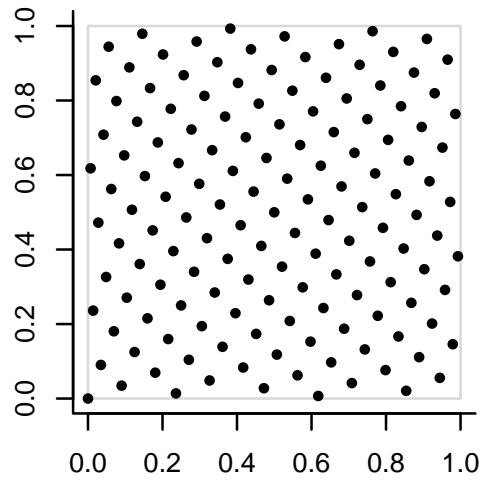
$$L^\perp = \{\mathbf{h} \in \mathbb{Z}^d \mid \mathbf{h}^\top \mathbf{Z} = 0 \pmod{n}\}$$

Smooth $f \implies \hat{f}$ decays with $|\hat{f}(\mathbf{h})|$.

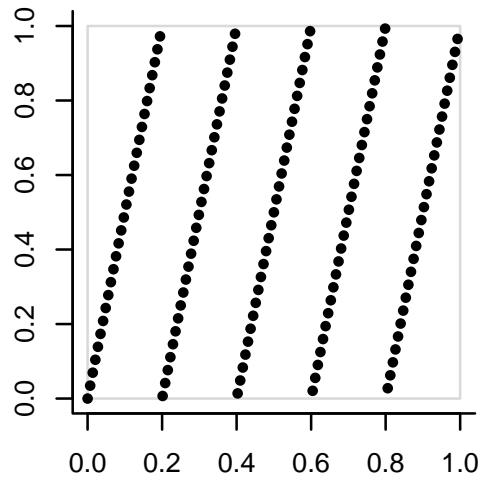
So search for \mathbf{Z} with $L^\perp \setminus \{\mathbf{0}\}$ 'far from origin'.

Dual lattice

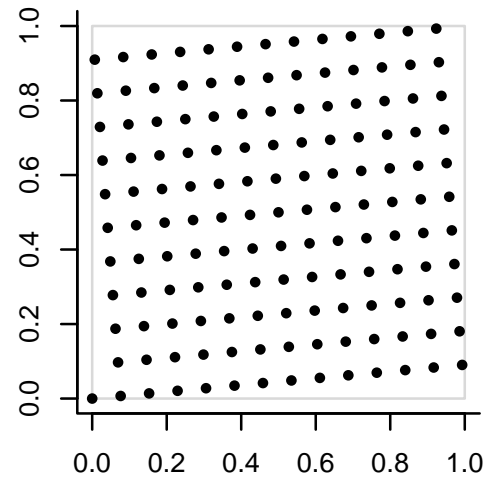
Some integration lattices



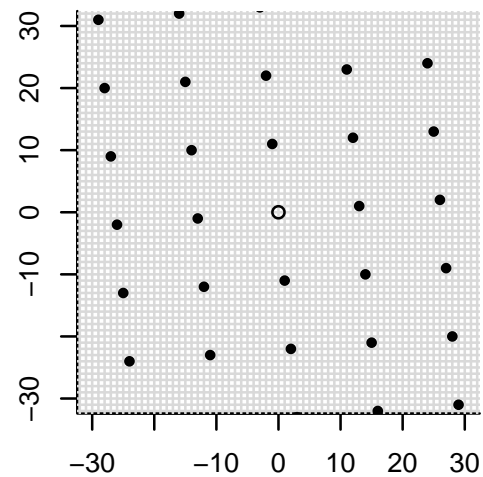
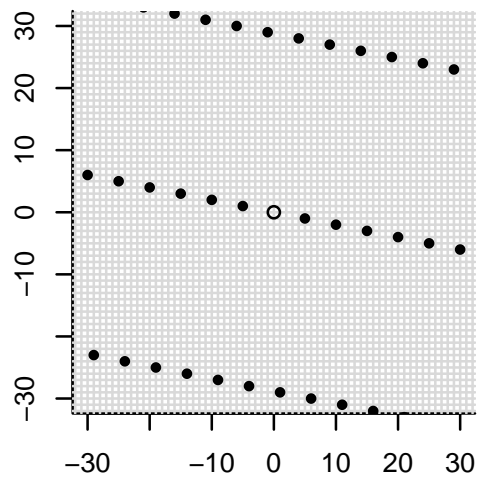
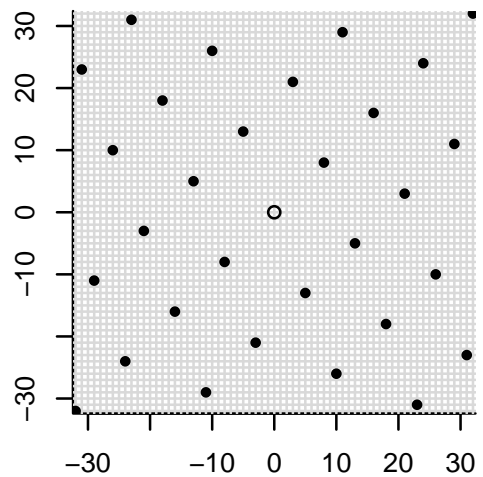
$n=144$ $z = (1,89)$



$n=144$ $z = (1,5)$



$n=144$ $z = (1,131)$



and their dual lattices

Periodizing transformation

$$\int_{[0,1]^d} f(\mathbf{x}) \, d\mathbf{x} = \int_{[0,1]^d} f(\tau(\mathbf{x})) J(\mathbf{x}) \, d\mathbf{x}$$

τ = transformation

J = Jacobian

Pick τ so that $J(\mathbf{x})$ vanishes on the boundary of $[0, 1]^d$.

Better: make r derivatives of $f \circ \tau \times J$ vanish on $\partial[0, 1]^d$.

Upshot

Good news: get very good convergence rate

Bad news: get strong curse of dimension in the constant.

QMC error estimation

$$|\hat{\mu} - \mu| \leq D_n^* \times V_{\text{HK}}(f)$$

Not a 100% confidence interval

- D_n^* is hard to compute
- V_{HK} harder to get than μ
- $V_{\text{HK}} = \infty$ is common, e.g., $f(x_1, x_2) = 1_{x_1+x_2 \leq 1}$
- We either get $|\hat{\mu} - \mu| < \infty$ or $|\hat{\mu} - \mu| \leq \infty$
(and maybe we already knew)

Also

Koksma-Hlawka is worst case. It can be very conservative.

Randomized QMC

- 1) Make $\mathbf{x}_i \sim \mathbf{U}[0, 1)^d$ individually,
- 2) keeping $D_n^* = O(n^{-1+\epsilon})$ collectively.

R independent replicates

$$\hat{\mu} = \frac{1}{R} \sum_{r=1}^R \hat{\mu}_r$$

$$\widehat{\text{Var}}(\hat{\mu}) = \frac{1}{R(R-1)} \sum_{r=1}^R (\hat{\mu}_r - \hat{\mu})^2$$

If $V_{\text{HK}}(f) < \infty$ then

$$\mathbb{E}((\hat{\mu} - \mu)^2) = O(n^{-2+\epsilon})$$

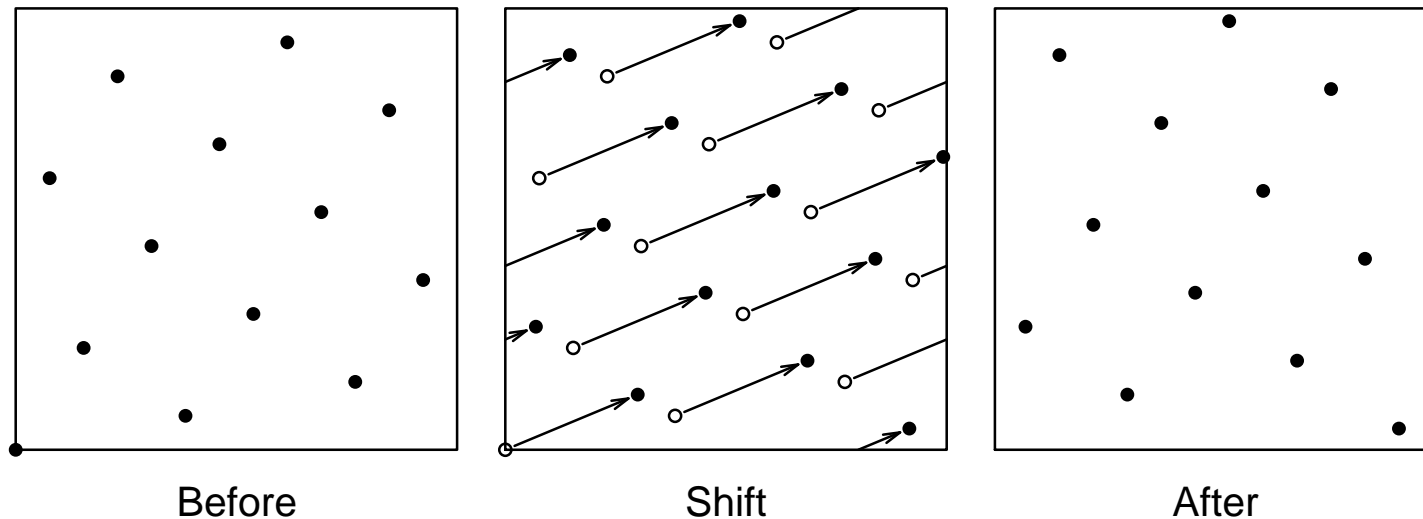
Random shift Cranley & Patterson (1976)

Scrambled nets O (1995,1997,1998)

Survey in L'Ecuyer & Lemieux (2005)

Rotation modulo 1

Cranley–Patterson rotation



Shift the points by $\mathbf{u} \sim \mathbf{U}[0, 1)^s$ with wraparound:

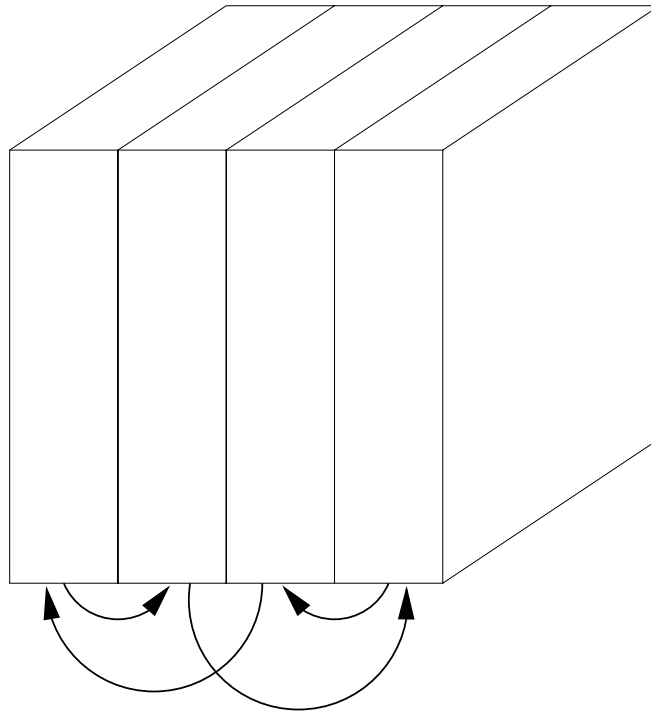
$$\mathbf{x}_i \rightarrow \mathbf{x}_i + \mathbf{u} \pmod{1}.$$

Commonly used on lattice rules.

Can also be used with nets.

At least it removes $\mathbf{x}_1 = 0$.

Digit scrambling



- 1) Chop the space into b slabs. Shuffle them.
- 2) Do the same within each of those b slabs.
- 3) And so on within b^2 , b^3 , \dots sub-slabs.
- 4) And the same for all s coordinates.

This operation yields $\mathbf{x}_i \sim \mathbf{U}[0, 1)^s$ and preserves the net property. [O \(1995\)](#)

Digit scrambling

Mathematically it is a permutation operation on the base b digits of

$$x_{ij} = \sum_{k=1}^K a_{ijk} b^{-k} \quad \rightarrow \quad \tilde{x}_{ij} = \sum_{k=1}^K \tilde{a}_{ijk} b^{-k}$$

The mapping $a_{ijk} \rightarrow \tilde{a}_{ijk} = \pi_{jk}(a_{ijk})$ preserves the net property.

In the previous 'nested' scramble π_{jk} also depends on $a_{ij1}, \dots, a_{ij,k-1}$

Simpler scrambles

Random linear permutations $a \rightarrow g + h \times a \pmod{p}$ (prime $p = b$)

$g \sim \mathbf{U}\{0, 1, \dots, b-1\}$, $h \sim \mathbf{U}\{1, 2, \dots, b-1\}$

Matousek (1998) has nested linear permutations

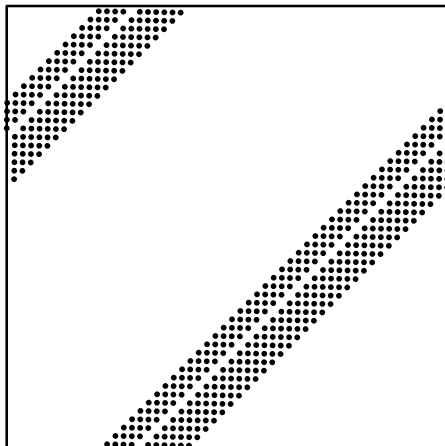
Digital shift $a_{ijk} \rightarrow a_{ijk} + g_{ij} \pmod{p}$

For $b = 2$ $\mathbf{x}_i \rightarrow \tilde{\mathbf{x}}_i = \mathbf{x}_i \oplus \mathbf{u}$ bitwise XOR

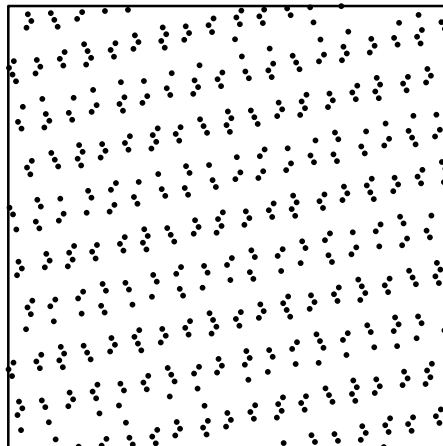
Example scrambles

Two components of the first 530 points of a Faure $(0, 53)$ -net in base 53.

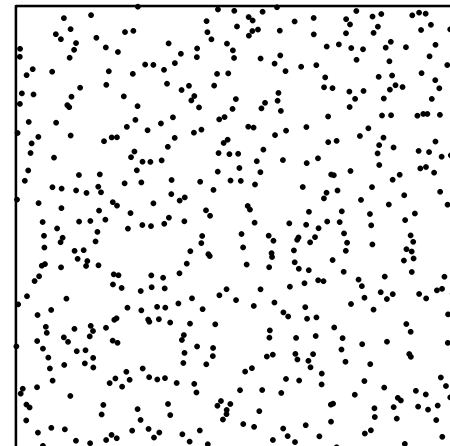
Randomized Faure points



Digital shift



Random linear



Nested uniform

The digital shift is much like a Cranley-Patterson rotation.

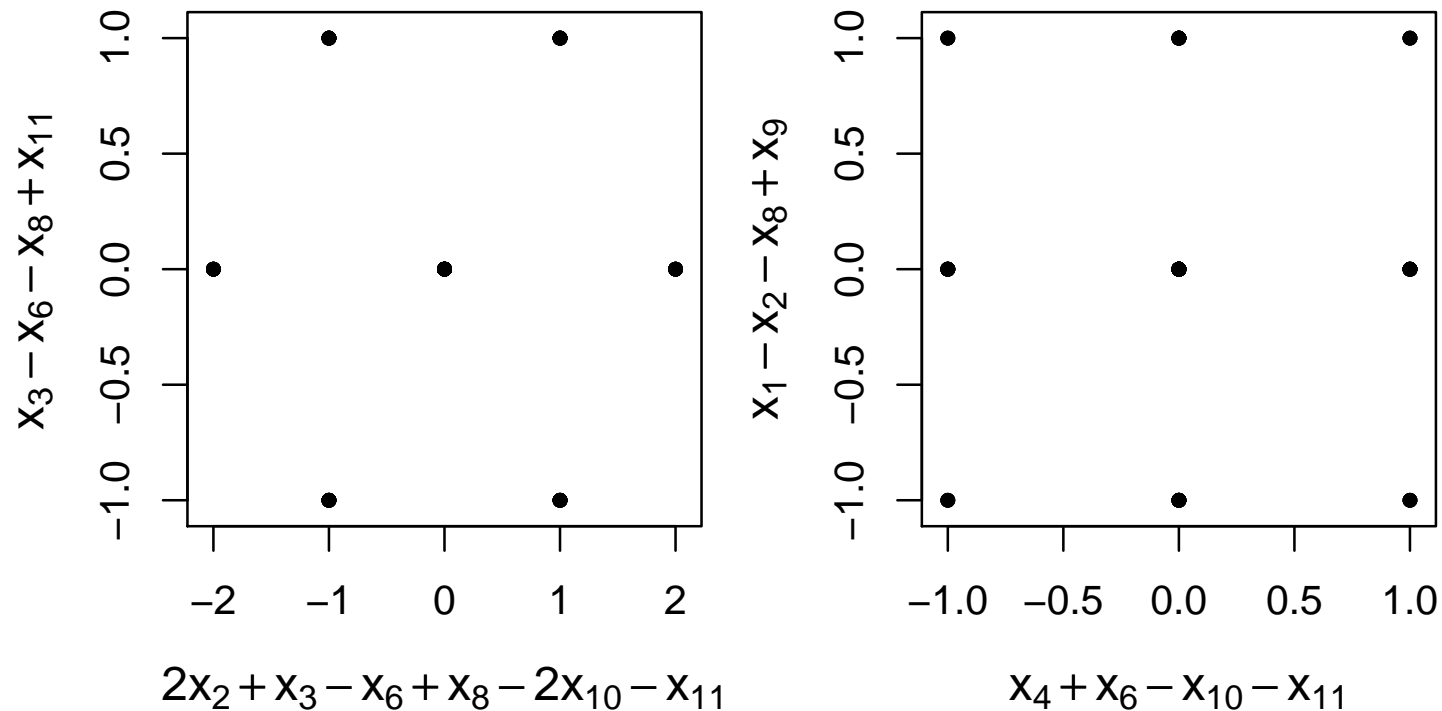
It uses just one random \mathbf{u} for all points: $\tilde{\mathbf{x}}_i = \mathbf{x}_i \oplus \mathbf{u}$.

Random linear [Matousek \(1998\)](#) and nested uniform [O \(1995\)](#) yield the same variance.

Unscrambled Faure

First $n = 11^2 = 121$ points of Faure $(0, 11)$ -net in $[0, 1]^{11}$.

Two projections of 121 Faure points



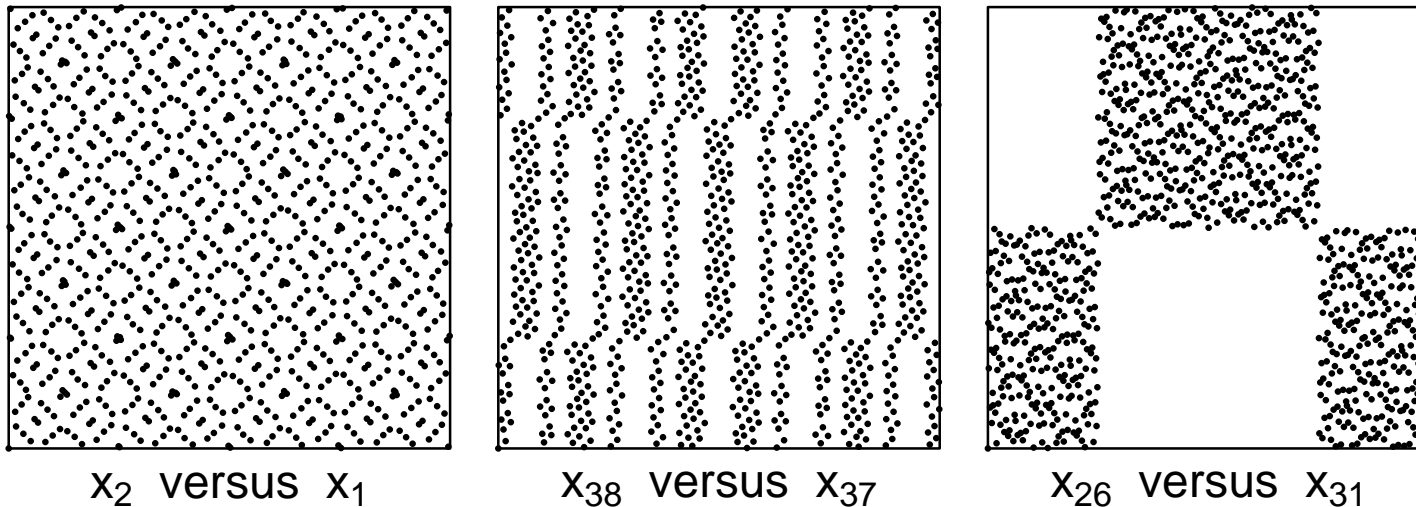
Unscrambled points are very structured.

Sobol' projections

Using code from [Bratley & Fox \(1988\)](#).

Implementations differ based on choices of 'direction numbers'.

Three projections of 1024 Sobol' points



At some higher sample size the empty blocks in x_{26} vs x_{31} fill in.

Expect empty blocks in higher order margins to last longer than those in low dimensional margins.

Scrambled net properties

Using $\sigma^2 = \int (f(\mathbf{x}) - \mu)^2 d\mathbf{x}$

If	Then	N.B.
$f \in L^2$	$\text{Var}(\hat{\mu}) = o(1/n)$	even if $V_{\text{HK}}(f) = \infty$
$f \in L^2$	$\text{Var}(\hat{\mu}) \leq \Gamma_{t,b,s} \sigma^2 / n$	for $t = 0$, $\Gamma \leq \exp(1) \doteq 2.718$
$\partial^{1,2,\dots,s} f \in L^2$	$\text{Var}(\hat{\mu}) = O(\log(n)^{s-1} / n^3)$	O(1997,2008)

$\Gamma < \infty$ rules out $(\log n)^{s-1}$ catastrophe at finite n .

Loh (2003) has a CLT for $t = 0$ (and **fully** scrambled points).

Geometrically

Scrambling breaks up the horizontal striping of the Faure nets.

Scrambling Sobol' points moves the full/empty blocks around.

Improved rate

RMSE is $O(n^{-1/2})$ better than QMC rate (cancellation).

Holds for nested uniform and nested linear scrambles.

I suspect that the CLT will **not** hold for random linear scrambles.

Hickernell and Yue found fourth moment differences in discrepancies for the two different kinds of scrambles.

I have a paper on recycling physical random numbers EJS (2009) where the result always has an asymptotically symmetric distribution which is never Gaussian. That is what lead me to guess (still unverified) that the CLT might not hold for the nested linear scramble of Matousek.

Scrambling vs shifting

Consider $n = 2^m$ points in $[0, 1)$.

QMC

van der Corput points $(i - 1)/n$ for $i = 1, \dots, n$.



Shift

Shift all points by $U \sim \mathbf{U}(0, 1)$ with wraparound.

Get one point in each $[(i - 1)/n, i/n)$



Scramble

Get a stratified sample, **independent** $x_i \sim \mathbf{U}[(i - 1)/n, i/n)$



Random errors cancel yielding an $O(n^{-1/2})$ improvement.

Higher order nets

Results from Dick, Baldeaux

Start with a net in $2s$ dimensions.

'Interleave' digits of two variables to make a new one

$$0.g_1g_2g_3 \cdots \text{ and } 0.h_1h_2h_3 \cdots \rightarrow 0.g_1h_2g_2h_2 \cdots .$$

Get a 'higher order' net in s dimensions.

Error is $O(n^{-2+\epsilon})$ under increased smoothness.

(partial twice wrt each input)

Scrambling gets RMSE $O(n^{-2-1/2+\epsilon})$

Even higher

Start with ks dimensions interleave down to s .

Get $O(n^{-k+\epsilon})$ and $O(n^{-k-1/2+\epsilon})$ (under still higher smoothness)

Upshot

Not yet widely used. Takes a lot of smoothness and a lot of inputs.

The curse of dimension

Curse of dimension: larger d makes integration harder.

$$C_M^r = \left\{ f : [0, 1]^d \rightarrow \mathbb{R} \mid \left| \prod_j \frac{\partial^{\alpha_j}}{\partial x_j^{\alpha_j}} f \right| \leq M, \sum_j \alpha_j = r, \alpha_j \geq 0 \right\}$$

Bahkvalov I:

For any $\mathbf{x}_1, \dots, \mathbf{x}_n \in [0, 1]^d$ there is $f \in C_M^r$ with $|\hat{\mu}_n - \mu| \geq kn^{-r/d}$

Ordinary QMC like $r = d$

Bahkvalov II:

Random points can't beat RMSE $O(n^{-r/d-1/2})$

Ordinary MC like $r = 0$

What if we beat those rates?

Sometimes we get high accuracy for large d .

It does not mean we beat the curse of dimensionality.

Bahkvalov never promised universal failure.

Only the existence of hard cases.

We may have just had an easy, non-worst case function.

Two kinds of easy

- Truncation: only the first $s \ll d$ components of x matter
- Superposition: the components only matter “ s at a time”

Either way

f might not be “fully d -dimensional”.

Lifted curse

Two main tools to describe it

- Weighted spaces and tractability
- ANOVA and effective dimension

Implications

Neither causes the curse to be lifted. They describe the happy circumstance where the curse did not apply.

Both leave important gaps described below.

Weighted spaces

Hickernell (1996), Sloan & Woźniakowski (1998)

$$\partial^u \equiv \prod_{j \in u} \frac{\partial}{\partial x_j} \quad \text{assume } \partial^{1:d} f \text{ exists}$$

Inner product, weights $\gamma_u > 0$

$$\langle f, g \rangle = \sum_{u \subseteq 1:d} \frac{1}{\gamma_u} \int_{[0,1]^u} \left(\int_{[0,1]^{-u}} \partial^u f(\mathbf{x}) d\mathbf{x}_{-u} \right) \left(\int_{[0,1]^{-u}} \partial^u g(\mathbf{x}) d\mathbf{x}_{-u} \right) d\mathbf{x}_u$$

Function ball $B_{\gamma,C} = \{f \mid \langle f, f \rangle_{\gamma} \leq C\}$

Small $\gamma_u \implies$ small $\|\partial^u f\|$ in ball.

Product weights

$\gamma_u = \prod_{j \in u} \gamma_j$ where γ_j decrease rapidly with j .

Now $f \in B_{\gamma,C}$ implies $\partial^u f$ small when $|u|$ large.

Tractability

$\text{Err}(d, n) \equiv$ Worst case error in function ball B_γ

$n_*(d, \epsilon) \equiv$ First n with $\text{Err}(d, n) \leq \epsilon \times \text{Err}(d, 0)$

Weak tractability

$$n_*(d, \epsilon) = \text{poly}(d, 1/\epsilon)$$

$$\limsup_{d \rightarrow \infty} \frac{\sum_{j=1}^d \gamma_j}{\log d} < \infty \quad \text{suffices}$$

Strong tractability

$$n_*(d, \epsilon) = \text{poly}(1/\epsilon)$$

$$\sum_{j=1}^{\infty} \gamma_j < \infty \quad \text{suffices}$$

Many versions and many contributors

Conditions above are for spaces described by [Sloan & Woźniakowski \(1998\)](#)

[Novak](#), [Wasilkowski](#), [Hickernell](#), [Dick](#), [Kuo](#) . . .

ANOVA and scrambled nets

Write the ANOVA of f

$$f(\mathbf{x}) = \sum_{u \subseteq \{1,2,\dots,s\}} f_u(\mathbf{x})$$

Term by term

$$\hat{\mu} = \sum_{u \subseteq \{1,2,\dots,s\}} \bar{f}_u \quad \text{where} \quad \bar{f}_u = \frac{1}{n} \sum_{i=1}^n f_u(\mathbf{x}_i)$$

$$\text{Var}(\hat{\mu}) = \sum_{u \subseteq \{1,2,\dots,s\}} \text{Var}(\bar{f}_u) \quad \text{O (1997)}$$

Expect $O(n^{-3+\epsilon})$ from small $|u|$ and $O(n^{-1})$ from large $|u|$

Lower order f_u are formed by integration. Can be smoother than f .

Griebel, Kuo, Sloan (2010, 2013, 2014)

Effective dimension $s \leq d$

Caflisch, Morokoff & O (1997)

Often f is dominated by its low order interactions.

Then RQMC may make a huge improvement.

Let $\sigma_u^2 = \text{Var}(f_u)$ variance component

Truncation sense

$$\sum_{u \subseteq 1:s} \sigma_u^2 \geq 0.99 \sum_{u \subseteq 1:d} \sigma_u^2$$

Superposition sense

$$\sum_{|u| \leq s} \sigma_u^2 \geq 0.99 \sum_{u \subseteq 1:d} \sigma_u^2$$

Mean dimension

$$\frac{\sum_u |u| \sigma_u^2}{\sum_u \sigma_u^2}$$

Liu & O (2006) (Easier to estimate \dots another story)

Example

Kuo, Schwab, Sloan (2012) consider quadrature for

$$f(\mathbf{x}) = \frac{1}{1 + \sum_{j=1}^d x_j^\alpha / j!}, \quad 0 < \alpha \leq 1.$$

For $\alpha = 1$ and $d = 500$

$R = 50$ replicated estimates of $\sum_v |v| \sigma_v^2 / \sigma^2$ using $n = 10,000$ had mean 1.0052 and standard deviation 0.0058.

Upshot

$f(\mathbf{x})$ is nearly additive

mean dimension between 1.00356 and 1.00684

(± 2 standard errors)

Weighted spaces

- 😊 γ_u ensure good performance for a class of f .
- 😞 Tractability conditions greatly restrict f .
- 😊 QMC can be tuned to given γ e.g., [Dick & Pillichshammer \(2010\)](#)
- 😞 We don't measure/observe a true γ .

Effective dimension

- 😊 Can be measured for a given f .
- 😞 Does not imply smooth interactions f_u .
- 😊 Anova components can be smoother than f [Sloan, Kuo, Griebel](#).

Description vs engineering

Effective dimension and weighted spaces both **describe** settings where f can be integrated well.

Some authors try to **engineer** changes in f to make it more suited to QMC.

E.g., cram importance into first few components of x .

MC vs QMC

MC places lots of effort on variance reduction.

For QMC we gain by reducing effective dimension.

Or Hardy-Krause variation (but there asymptotics are slow).

E.g., turning $u \sim \mathbf{U}[0, 1]^d$ into $z \sim \mathcal{N}(0, \Sigma)$

The choice of $\Sigma^{1/2}$ affects QMC performance.

Caflich, Morokoff & O (1997)

Acworth, Broadie & Glasserman (1998)

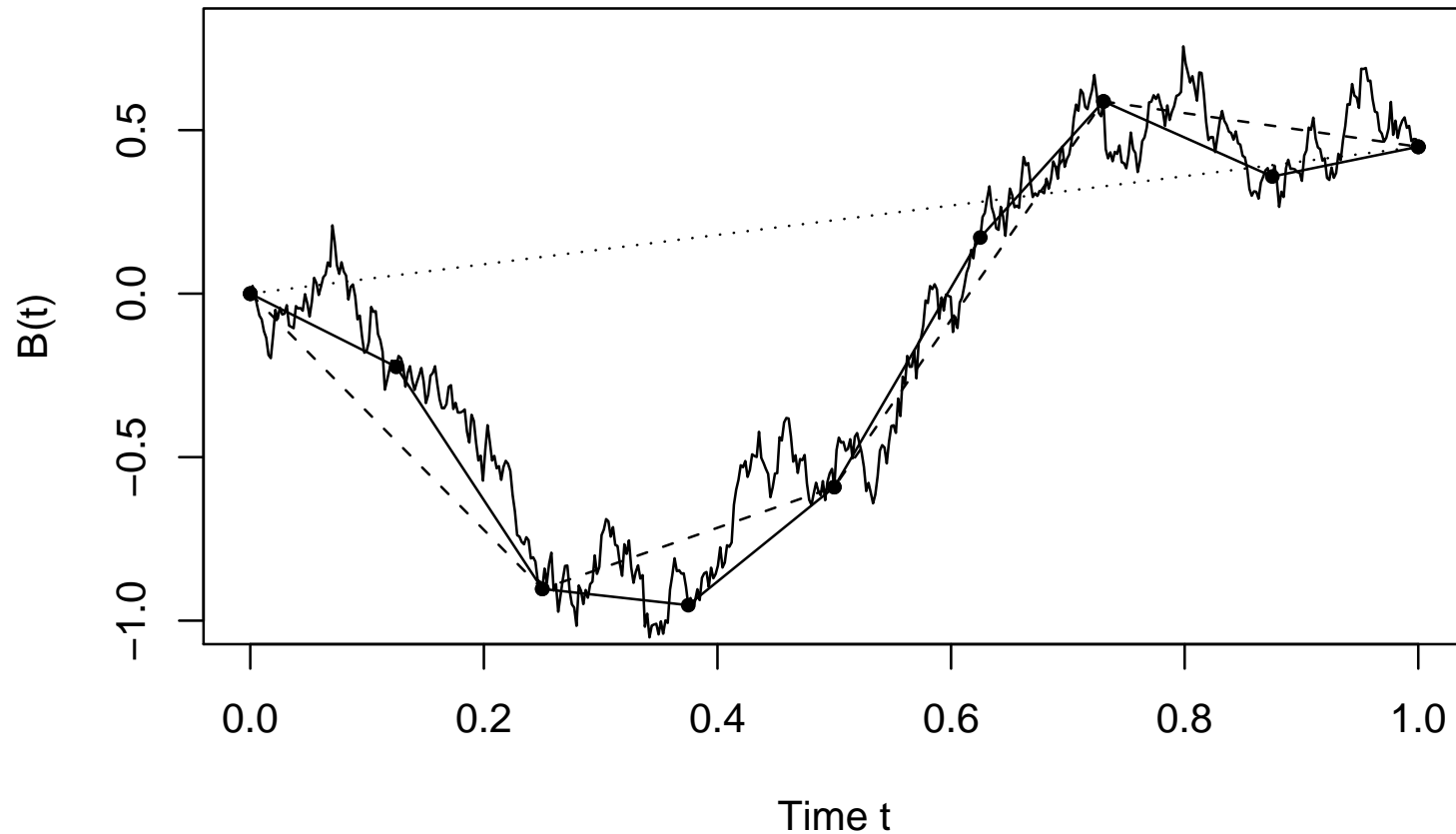
Imai & Tan (2014)

Sampling Brownian motion

Feynman-Kac/Brownian bridge

First few variables define a 'skeleton'. The rest fill in.

Brownian bridge construction of Brownian motion



See also [Mike Giles++](#) on multi-level MC.

MCMSki, January 2016

Acworth, Broadie and Glasserman simulate BM in terms of its principal components. If the first few (or any few) PCs dominate the evaluation function, then it reduces effective dimension.

There can be no universal BM generator that is best for all integrands. Papageorgiou (2002) exhibits a financial option for which the ordinary step by step construction of Brownian motion works better than the Brownian bridge construction.

Imai and Tan have a series of papers tuning the generation of Brownian motion to specific integrands.

The book edited by Schreider (1964) “Method of statistical testing, Monte Carlo method” includes a Chapter II by I. M. Sobol’ in which the Brownian bridge construction is described. Sobol’ cites a 1958 article by Gel’fand, Frolov and Chentsov “The evaluation of Wiener integrals by the Monte Carlo method”. Izv. vys. uch. zav. ser. matem., No. 5, 32–45

Choosing γ

Each γ corresponds to a reproducing kernel Hilbert space (RKHS)

The question

Which RKHS should we use in a given problem?

\mathcal{H}_1 or \mathcal{H}_2 or \dots or \mathcal{H}_J \dots

- 1) sometimes $f \in \mathcal{H}_j$ **all** $j = 1, \dots, J$
and $f \in \mathcal{H}_1$ vs \mathcal{H}_2 have very different implications
- 2) sometimes f belongs to **none** of them.
while $|f - \tilde{f}| \leq \epsilon$ where $\tilde{f} \in \mathcal{H}$

Bayes and empirical Bayes ideas might help choose \mathcal{H}

Maybe we want an \mathcal{H} where f is ‘typical’.

A natural γ has

$$\gamma_u \propto \int_{[0,1]^d} (\partial^u f(\mathbf{x}))^2 d\mathbf{x}$$

NB: the constant of proportionality is also important.

Asymptotics

Two ways to think about asymptotics.

View 1

We really want to know what happens in the limit as $n \rightarrow \infty$.

View 2

We want an approximation for finite n , like $\frac{1}{n} \approx 0$.

For QMC

If the dimension is large, then the asymptotics solve view 1 not view 2.

Practical n are likely to be on an 'initial transient'

Sometimes the asymptote 'sets in' at $n \approx 4^d$.

For MC

$\text{Var}(\hat{\mu}) = \sigma^2/n$ 'sets in' for $n \geq 1$.

Sequential QMC

JRSS-B discussion paper by [Chopin & Gerber \(2015\)](#)

N particles in \mathbb{R}^d for T time steps

Advance \mathbf{x}_{nt} to $\mathbf{x}_{n,t+1} = \phi(\mathbf{x}_{nt}, \mathbf{u}_{n,t+1})$ each $\mathbf{u} \in [0, 1)^S$.

Do them all with $U_{t+1} \in [0, 1)^{N \times S}$. Usually IID.

For QMC

We want rows of U_{t+1} to be ‘balanced’ wrt position of \mathbf{x}_t

Easy if $d = S = 1$ [Lécot & Tuffin \(2004\)](#)

Take $U_{t+1} \in [0, 1)^{N \times 2}$

Sort first column along $\mathbf{x}_{n,t}$. Update points with remaining column.

What to do if $d > 1$?

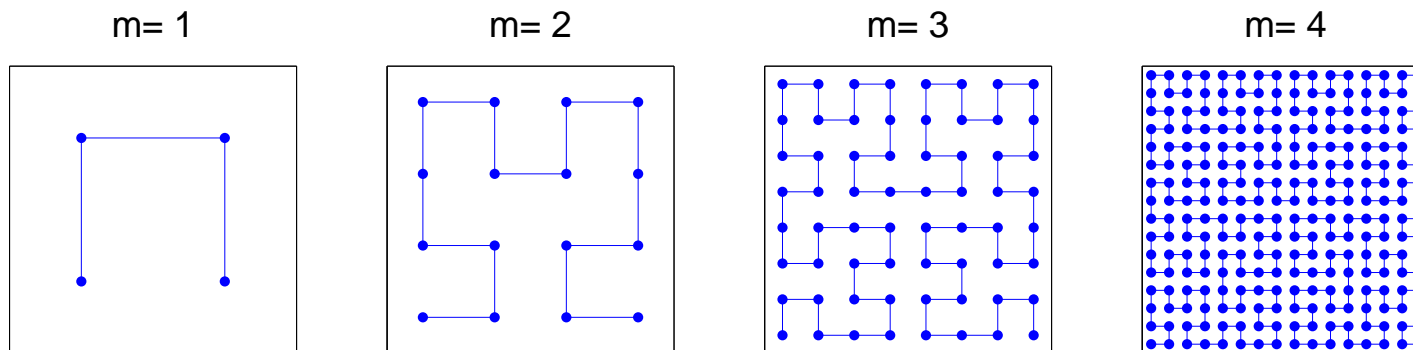
Alignment is tricky because \mathbb{R}^d is not ordered.

Related work by [L'Ecuyer, Lécot, L'Archevêque-Gaudet](#) on array-RQMC.

d -dimensional alignments. Huge variance reductions. Limited theory.

Using Hilbert curves

These are space-filling curves as $m \rightarrow \infty$:



Chopin & Gerber thread a curve through $\mathbf{x}_{n,t} \in \mathbb{R}^d$

Align particles via first column of RQMC points $U_{t+1} \in [0, 1)^{N \times (1+S)}$.

Use remaining S columns to advance them.

Results

Squared error is $o(N^{-1})$.

Good empirical results for modest d

NB:

I left out lots of details about particle weighting.

Recipe for QMC in MCMC

For MCMC we follow one particle.

- 1) Step $\mathbf{x}_i \leftarrow \psi(\mathbf{x}_{i-1}, \mathbf{v}_i)$ for $\mathbf{v}_i \in (0, 1)^d$.
- 2) For $\mathbf{v}_1 = (u_1, \dots, u_d)$, $\mathbf{v}_2 = (u_{d+1}, \dots, u_{2d})$ etc.
 n steps require $u_1, \dots, u_{nd} \in (0, 1)$
- 3) MCMC uses $u_i \sim \mathbf{U}(0, 1)$
- 4) Replace IID by $N = nd$ 'balanced' points

Reasons for caution

- 1) We're using 1 point in $[0, 1]^{nd}$ with $n \rightarrow \infty$
- 2) Our \mathbf{x}_i won't be Markovian

QMC \cap MCMC

Early references

Chentsov (1967)

Plugs in 'completely uniformly distributed' points. (defined later)

Samples in finite state space by inversion.

Shows consistency.

Beautiful coupling argument. (Resembles coupling from the past.)

Sobol' (1974)

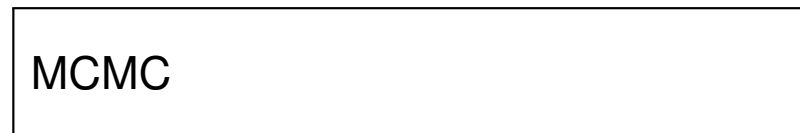
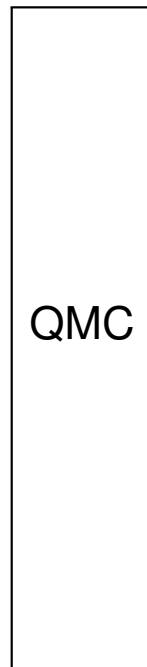
Has $n \times \infty$ points $x_{ij} \in [0, 1]$

Samples from a row until a return to start state, then goes to next row

Gets rate $O(1/n) \dots$ if transition probabilities are $a/2^b$ for integers a, b

$$\text{MCMC} \approx \text{QMC}^T$$

Method	Rows	Columns	
QMC	n points	d variables	$1 \leq d \ll n \rightarrow \infty$
MCMC	r replicates	n steps	$1 \leq r \ll n \rightarrow \infty$



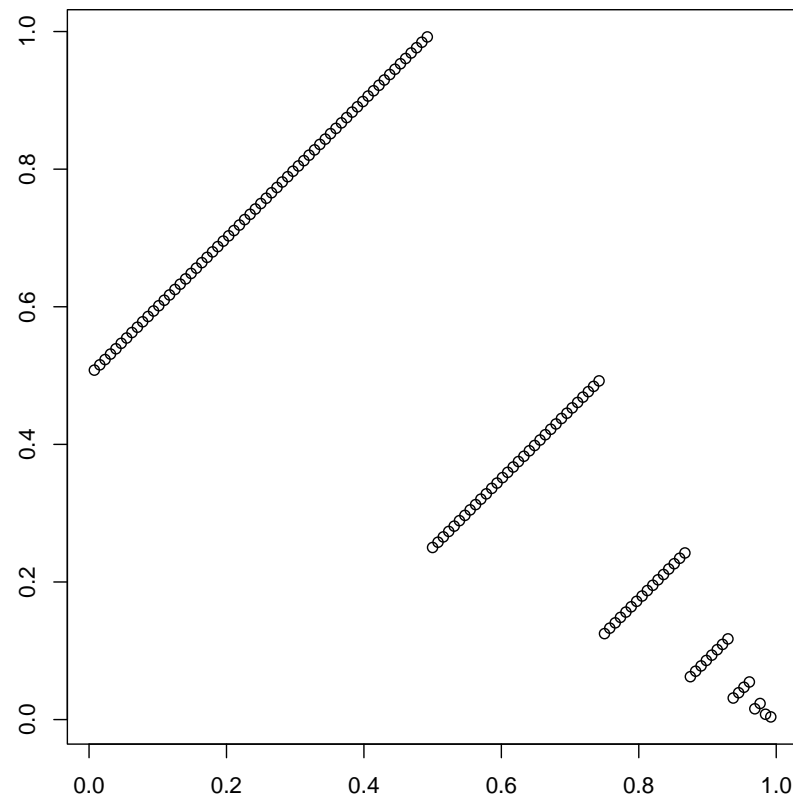
QMC based on equidistribution

MCMC based on ergodicity

Severe failure is possible

van der Corput $u_i \in [0, 1/2) \iff u_{i+1} \in [1/2, 1)$

u_{i+1} VS u_i



High proposal \iff low acceptance and vice versa

Morokoff and Caflisch (1993) describe heat particle leaving region

Completely uniformly distributed

$u_1, u_2, \dots \in [0, 1]$ are CUD if

$D_n^*(z_1, \dots, z_n) \rightarrow 0$, where

$z_i = (u_i, \dots, u_{i+d-1})$

For all $d \geq 1$

Overlapping blocks

$$z_1 = (u_1, \dots, u_d)$$

$$z_2 = (u_2, \dots, u_{d+1})$$

$$\vdots \quad \vdots$$

$$z_n = (u_n, \dots, u_{n+d-1})$$

Chentsov (1967) shows we can use non-overlapping blocks

$$v_i = (u_{d(i-1)+1}, \dots, u_{di}) \quad \forall d$$

CUD ctd

CUD \equiv one of Knuth's **definitions** of randomness

Recommendations

- 1) Use all the d -tuples from your RNG
- 2) Be sure to pick a small RNG

As considered in

Niederreiter (1986)

Entacher, Hellekalek, and L'Ecuyer (1999)

L'Ecuyer and Lemieux (1999)

Some results

- 1) Metropolis on CUD points consistent in finite state space [O & Tribble \(2005\)](#)
- 2) Weakly CUD random points ok [Tribble & O \(2008\)](#)
- 3) Consistency in continuous problems (Gibbs and Metropolis)
[Chen, Dick & O \(2010\)](#)
- 4) [Chen \(2011\)](#) rate $O(n^{-1+\epsilon})$ for some smooth ARMA

Related refs

Liao (1998)	reorders QMC points for MCMC
Craiu & Lemieux (2007)	QMC in multiple-try Metropolis
Lemieux & Sidorsky (2006)	QMC in exact sampling
Lemieux, Ormoneit, Fleet (2001)	QMC for particle filters
Frigessi, Gäsemyr, Rue (2000)	MCMC \cap antithetics
Craiu, Meng (2004)	MCMC \cap Latin hypercubes
Propp (2004)	Rotor-Router

Thesis of Tribble (2007)

Results on CUD and weak CUD

Constructions of small RNGs

Examples with Gibbs and Metropolis

Variance reductions from Tribble

Data sets	$n = 2^{10}$		$n = 2^{12}$		$n = 2^{14}$	
	min	max	min	max	min	max
Pumps ($d = 11$)	286	1543	304	5003	1186	16089
Vasorestriction ($d = 42$)	14	15	56	76	108	124

Pumps: hierarchical Poisson-Gamma model.

Vasorestriction: probit model 3 coefficients, 39 latent variables.

Min & max variance reductions for all pump and all non-latent vaso. parameters.

Details

Targets are posterior means of parameters.

CUD points were LFSR, with Cranley-Patterson rotations.

Chen, Dick & O Annals Stat.

Chen, Dick & O extend consistency to continuous state spaces.

MCMC remains consistent when driven by u_1, u_2, \dots , if

- 1) u_i are CUD (or CUD in probability)
- 2) m -step transitions are **Riemann** integrable $\forall m \geq 1$, and
- 3)
 - for Metropolis-Hastings: there is a **coupling** region
(Independence sampler can have one)
 - for Gibbs: there is a **contraction** property
(Gibbs for probit model proven to contract)

Either way, the chain has to forget its past.

That could be by regeneration or exponential decay.

Aside: Henri Lebesgue made our lives so much easier!

Small RNGs

Matsumoto & Nishimura sent us some small linear feedback shift register RNGs

Similar equidistribution properties to “small Mersenne twisters” but not necessarily the same constructions.

They come in sizes $M = 2^m - 1$ for $10 \leq m \leq 32$.

$$u_1, u_2, \dots, u_M$$

We explore them for some simulations.

Prepend one or more 0s:

$$0, \dots, 0, u_1, \dots, u_M$$

put into a matrix and apply Cranley-Patterson rotations

Summary

Bivariate Gaussian	apparent better convergence rate for mean
Hit and run, volume estimator	no improvement
M/M/1 queue, average wait	mixed results
Garch	some big improvements
Heston stochastic volatility	big improvements for in the money case

Synopsis

The smoother the problem, the more CUD points can improve.

Same as for finite dimensional QMC.

What next?

QMC helps MCMC the most when transitions are smooth.
i.e., no acceptance rejection.

Some thoughts

- Gibbs is promising. Maybe tweak JAGS.
- Maybe Hamiltonian dynamics using high acceptance.
- When there is acceptance-rejection, maybe the optimal acceptance rate for QMC is higher.

Thanks

- QMC \cap MCMC co-authors:
Seth Tribble, Su Chen, Josef Dick, Makoto Matsumoto, Takuji Nishimura
- MCMSki organizers
- Nicolas Chopin
- NSF