

Monte Carlo: Techniques and Theory

Ronald Kleiss¹
IMAPP, Radboud University of Nijmegen

version of February 7, 2019



John von Neumann
the other godfather



Stanislaw Ulam
the other godfather



Nick Metropolis
the other godfather

*‘We guarantee that each number
is random
individually, but we don’t
guarantee that more than one of them
is random’*

Unnamed programming consultant quoted in [1]

¹R.Kleiss@science.ru.nl

Contents

1	Introduction: randomness and probability	11
1.1	Protohistory	11
1.2	Random number streamsRandom numbers!stream	11
1.2.1	Random numbers: <i>circulus in probando</i>	11
1.2.2	What is a random stream? Relying on probability	11
1.2.3	What is probability? Relying on a random stream	12
1.2.4	A difference between physics and mathematics	12
1.3	Miscellaneous probability items	12
1.3.1	Some notation in these notes	12
1.3.2	Moments and characteristic function	13
1.3.3	The Chebyshev-Bienaymé inequality	13
1.3.4	The Central Limit Theorem	14
2	Monte Carlo integration	16
2.1	The Monte Carlo idea	16
2.1.1	Point sets and expectation values	16
2.1.2	Integration as archetype	16
2.1.3	Point sets, ensembles, and the Leap Of Faith	17
2.2	Finding estimators	17
2.2.1	Direct sums, unequal-sums, and expectation values	17
2.2.2	The Monte Carlo estimators	18
2.2.3	Positivity of E_2 and E_4	19
2.3	Estimating in practice	20
2.3.1	Improved estimators	20
2.3.2	Numerical stability; the extended CGL algorithm	20
2.3.3	How to <i>do</i> Monte Carlo integration	21
2.3.4	How to <i>report</i> Monte Carlo integration	22
2.3.5	How to <i>interpret</i> Monte Carlo integration	22
2.4	A test case	23
2.5	Exercises	26
3	Random number generation	28
3.1	Introduction to random number sources	28
3.1.1	Natural <i>vs</i> Pseudo	28
3.2	Pseudo-random number streams	29
3.2.1	The structure of pseudo-random number algorithms	29

3.2.2	The set of all algorithms	30
3.2.3	The concept of an arbitrary algorithm	30
3.2.4	Lessons from random algorithms	31
3.2.5	Shift-register PRNGs	32
3.3	Bad and good algorithms	32
3.3.1	Bad: the midsquare algorithm	32
3.3.2	Bad: the chaos approach and the logistic map	34
3.3.3	Maybe not so bad: linear congruential methods	35
3.3.4	Horrible: RANDU, and a possible redemption	37
3.3.5	Good: RCARRY and RANLUX	37
3.3.6	Good: the Mersenne Twister	39
3.4	Exercises	41
4	Testing PRNGs	43
4.1	Empirical testing strategies and doubts	43
4.1.1	The Leeb Conundrum: too much of a good thing	43
4.1.2	Any test is a uniformity test	44
4.1.3	The χ^2 characteristic	44
4.2	Theoretical testing strategies	45
4.2.1	The number-to-number correlation	45
4.2.2	The spectral test	45
5	Quasi-Monte Carlo	46
5.1	Generalities of QMC	46
5.1.1	The New Leap of Faith	46
5.1.2	The mechanism of error improvement	47
5.2	Error estimators	48
5.2.1	The first-order estimate	48
5.2.2	The second-order estimate	48
5.2.3	Payback time: Lack of Leap of Faith is Punished	48
6	Nonuniformity of point sets	49
6.1	Measures of nonuniformity: Discrepancy	49
6.1.1	The star discrepancy	49
6.1.2	Random <i>vs</i> Regular: Translation <i>vs</i> Rotation	50
6.1.3	The Roth bound	51
6.1.4	The Koksma-Hlawka inequality	51
6.1.5	The Wiener measure and the Woźniakowski Lemma	51

6.2	Measures of nonuniformity: Diaphony	53
6.2.1	Fourier problem classes	53
6.2.2	Fourier diaphony	54
6.2.3	Choosing your strengths: examples of diaphony	55
6.3	QFT for diaphony	56
6.3.1	The distribution of diaphony	56
6.3.2	Feynman rules for diaphony in the large- N limit	57
6.3.3	Collecting bracelets	58
6.3.4	The diaphony distribution for large N	60
6.3.5	The saddle-point approximation	61
6.3.6	$1/N$ corrections to the diaphony distribution	62
6.3.7	The two-point function	63
6.3.8	Testing too much: the Dirac limit	65
6.4	Measures of nonuniformity: χ^2	66
6.4.1	The χ^2 as a discrepancy	66
6.4.2	Large- N results for χ^2	67
6.4.3	Two-point function and $1/N$ corrections for χ^2	68
7	Superuniform point sets	70
7.1	Fixed point sets <i>vs</i> streams	70
7.1.1	Diaphony minimisation	70
7.1.2	Korobov sequences: good lattice points	71
7.2	QRNG algorithms	72
7.2.1	Richtmeyer-Kronecker streams	72
7.2.2	Excursion into fractions (cont'd)	73
7.2.3	Rational approximations to irrationals	74
7.2.4	Almost-equidistancy for Richtmeyer sequences	75
7.2.5	van der Corput streams	76
7.2.6	Van der Corput sequences in more dimensions	78
7.2.7	Niederreiter streams	80
8	Variance reduction	81
8.1	Stratified sampling	81
8.1.1	General strategy	81
8.1.2	An example: VEGAS	81
8.1.3	An example: PARNI	81
8.2	Importance sampling	81
8.2.1	General strategy	81

8.2.2	Multichanneling	81
9	Non-uniform PRNGs	82
9.1	The Art of Transforming, Rejecting, and Being Smart	82
9.2	The UA formalism	82
9.2.1	Unitary algorithms as words and as pseudocode	82
9.2.2	Inversion of variates in UA	84
9.2.3	Rejection of variates in UA	85
9.3	Repertoire and the Rule of Nifty	87
9.3.1	Building up a repertoire	87
9.3.2	The normal distribution: the Box-Müller algorithm	88
9.3.3	The Euler algorithm	89
9.3.4	The Kinderman-Monahan algorithm	91
9.4	Random-walk algorithms	92
9.4.1	The Metropolis algorithm	93
9.4.2	An elementary case study for Metropolis	95
9.4.3	Applications of the Metropolis algorithm	96
9.4.4	Gibbs sampling	97
9.4.5	An elementary case study for Gibbs	98
10	Phase space algorithms for particle physics	100
10.1	The uniform phase space problem in particle phenomenology	100
10.2	Two-body phase space	100
10.2.1	The two-body algorithm	100
10.2.2	Two-body reduction	101
10.3	The relativistic problem	102
10.3.1	Two-body reduction algorithm	102
10.3.2	Massless RAMBO	103
10.3.3	Inclusion of masses	105
10.4	Nonrelativistic phase space: BOLTZ	107
11	Appendices	110
11.0.1	Falling powers	110
11.0.2	Relations between direct sums and unequal-sums	111
11.0.3	An exponential sum and the Poisson formula	111
11.0.4	About the integral (40)	112
11.0.5	Selfies	112
11.0.6	Serial correlation in a real-number model	113

11.0.7	The two-point function for the Euler diaphony	114
11.0.8	Rational denominators for continued fractions	114

List of Algorithms

1	Extended CGL algorithm for numerically safe updating of running Monte Carlo estimators	22
2	The RCARRY algorithm using a cyclic register.	38
3	The two major steps of the MT algorithm.	41
4	The van der Corput transform $\phi_b(n)$	77
5	The Box-Müller algorithm	89
6	The Euler density with parameters p_1, p_2, \dots, p_n	90
7	Generating the Cauchy density by ratio of uniforms	92
8	Two-body phase space with masses $m_{1,2}$ and total invariant energy $\sqrt{s} > m_1 + m_2$	101
9	Lorentz boost from P^μ , with $P^2 = s$, at rest to given form, applied on vector p^μ . The resultant vector is q^μ	101
10	The Platzner algorithm for $n \geq 3$	103
11	The RAMBO algorithm for n momenta with total invariant mass squared s	105
12	Lorentz boost from P^μ , with $P^2 = s$, to rest from given form, applied on vector p^μ . The resultant vector is q^μ	105
13	Giving masses to massless momenta	106
14	The BOLTZ algorithm for total energy U and masses $m_{1,2,\dots,n}$	108

The Buffon fragment

What follows is translated from the ‘Buffon fragment’, clay tablet in Akkadian discovered at the site of *Jebel-i-Qurul*, most likely the library of the temple precinct of the deity Nisaba, which contained a school of the *tupsar enuma Anu Enlil* scribes, who specialized in astronomy and astrology¹; believed to be based on an earlier Sumerian original, primarily because of the reference to *absu*. Several scholars², however, maintain that this fragment is a forgery.

The main text

Master: Inside the lowly reedstalks! thou mayest find, O my disciple, the secret [of] the temple’s column ; yea, verily, the secret of its girth to its width³, from the river’s reeds! To penetrate [the secret], to gain the column’s wisdom, thou shalt go to the water’s edge to gather reeds and bring them together, yea, even as many as thou canst gather; and [thou] shalt cut [them], so that none shall surpass the others, nor one be less than any⁴; and thou shalt also take clay

from the water’s edge, even as much as thou canst gather, [and] bring it to the scribe’s apprentice. And [the scribe’s apprentice] shall shape thereof [a tablet of] two cubits; and [the scribe’s apprentice] shall draw many [lines on the tablet] so that none approach nor separate; and the empty space⁵ between [the lines] shall be as one reed-stalk, so that it neither crosses [them], nor shall it fall short: but the reed will be like unto a bridge from one line to another⁶. Thou shalt empty thy mind of all [thought], [thou shalt] void thy spirit of all purpose; and thou shalt throw [the stalks] down onto the tablet, yea, and scatter them, even like unto chaff that is scattered by the wind [on the] threshing floor⁷. And the reeds that cross [a line], those thou shalt gather together in thy hands, but the reeds that do not cross thou shalt not [gather together]. And the multitude in the number of the reeds in thy hands⁸, thou shalt [take] anew⁹. And behold! it is as the column, yea, even as the girth of the column to its width [...]

Disciple: O my master, if I [perform] this task, and my brother [performs]

⁵literally, ‘absu’ i.e. the watery abyss

⁶ i.e. draw parallel lines at a distance of precisely one stalk’s length

⁷i.e. throw the reedstalks at random, with no preconceived pattern

⁸i.e. the the total number of stalks divided by the number of retained stalks

⁹i.e. multiply by two

¹ A.L. Oppenheim, *Ancient Mesopotamia* (univ. of Chicago Press, 1977), p242.

² see, for instance, von Däniken’s reference to the Book of Dzyan (1968).

³i.e. the value of π

⁴i.e. cut them to precisely the same length

this task, and all my brethren [perform] this task, shall [we] not then [approach] closer to the secret of the temple [’s column]?

Master: Verily, thou speakest [with] wisdom, O my [disciple]; for as a single stalk leadeth not towards knowledge, and giveth not the secret; so many [stalks] shall reveal much of [the secret]. Yet lo! the secret is revealed ever more slowly to the diligent [...]¹⁰

Analysis of the prescription

We assume that the stalks are straight lines of unit length; likewise the lines on the tablet are parallel straight lines with unit separation. Consider a reed-stalk that makes an angle ϕ with respect to the parallels. Its projection orthogonal to the parallels then has length $\sin(\phi)$. Supposing that the distribution of the stalks over the tablet is translationally invariant, this gives the probability that the given stalk will intersect one of the lines: note that this relies on the fact that the tablet be large enough to contain all the thrown stalks, hence the reference to ‘two cubits’¹¹. The procedure requires throwing the stalks with ran-

domly chosen orientation, that is, ϕ is a random variable distributed uniformly between 0 and π (note that $\phi \rightarrow \phi + \pi$ gives the same situation, except for a change in the stalk’s orientation which is irrelevant). The expected probability for a given stalk to cross one of the parallels is therefore

$$\langle \phi \rangle = \frac{1}{\pi} \int_0^\pi \sin(\phi) d\phi = \frac{2}{\pi} ,$$

so that we arrive at

$$\pi = \frac{2}{\langle \phi \rangle} .$$

The expected probability is measured by using many stalks and estimating $\langle \phi \rangle$ by the value of x , where

$$x = \frac{\text{no. of stalks crossing a line}}{\text{total number of stalks}} ,$$

which proves the validity of the algorithm. The estimate for $\langle \phi \rangle$ improves with increasing number of stalks, as suggested by the disciple’s question. On the other hand, as indicated by the master, the convergence to the exact answer is only asymptotic. The expected error in the estimate after N stalks have been thrown can be computed to be

$$\begin{aligned} |x - \langle \phi \rangle| &\approx \sqrt{\frac{\frac{2}{\pi} \left(1 - \frac{2}{\pi}\right)}{N}} \\ &\approx \frac{0.481}{\sqrt{N}} . \end{aligned}$$

¹⁰A reference to $1/\sqrt{N}$ convergence?

¹¹The mesopotamian cubit is about 51.86 cm, from the specimen discovered by E. Unger at Nippur (*Acta praehistorica et archaeologica* Vol 7. Berliner Gesellschaft für Anthropologie, Ethnologie und Urgeschichte, Hessling Verlag, 1976).

To obtain an accuracy of 2 decimal digits (*i.e.* to get $\pi \approx 3.14$) one would need about 23,000 stalks, which may explain ‘as many as thou canst gather’. To obtain the next digit would necessitate the use of 2.3 million stalks.

The fragment breaks off in the middle of the disciple’s exclamation: ”Woe is me! The reeds are a heavy [burden], a terrible multitude, an angry host [...]”

An early scientific application Monte Carlo

Lord Kelvin reports [2] on an early Monte Carlo method for simulating the motion of a particle in a volume with roughened edges, in order to examine the equidistribution of energy:

‘[...] I have evaded the difficulty in a manner thoroughly suitable for thermodynamic application such as the kinetic theory of gases. I arranged to draw lots for out of the 199 points dividing AB into 200 equal parts. This was done by taking 100 cards*, 0, 1 98, 99, to represent distances from the middle point, and, by the toss of a coin, determining on which side of the middle point it was to be (plus or minus for head or tail, frequently changed to avoid possibility of error by bias). The draw for one of the hundred numbers (0 99) was taken after very thorough shuffling of the cards in each case [...].’

The footnote reads:

* ‘I had tried numbered billets (small squares of paper) drawn from a bowl, but found this very unsatisfactory. The best mixing we could make in the bowl seemed to be quite insufficient to secure equal chances for all the billets. Full sized cards like ordinary playing-cards, well shuffled, seemed to give a very fairly equal chance to every card. Even with the full-sized cards, electric attraction sometimes intervenes and causes two of them to stick together. In using one’s fingers to mix dry billets of card, or of paper, in a bowl, very considerable disturbance may be expected from electrification.’

1 Introduction: randomness and probability

1.1 Protohistory

Monte Carlo methods are those numerical approaches to any problem in which *at least one* random number is used to obtain an outcome. The idea is not new, see for instance [2]. It came to fruition with the advent of computers in the 1940's, mainly under the influence of WW2 efforts [3, 4]. One may distinguish Monte Carlo *integration* (of functions) and Monte Carlo *simulation* (of processes). Formally these amount to the same.

1.2 Random number streamsRandom numbers!stream

1.2.1 Random numbers: *circulus in probando*

The concept of what constitutes a set of random numbers is surprisingly tenuous. *Any* given set of numbers can be subjected to exhaustive analysis and so be 'shown' to be not 'fortuitous' but 'determined'¹². Therefore, a 'true' set of random numbers should be considered rather as a *stream* of numbers, like a tap that can be turned on and, maybe very much later on, be turned off, or left running indefinitely. The idea is that it need never stop, and the collection of numbers can in principle grow without limit.

⌘ Notions of randomness as 'computational complexity' and 'Kolmogorov complexity' are essentially defined for finite sets.

1.2.2 What is a random stream? Relying on probability

The (to my mind) most operationally useful definition of 'truly random numbers' resides in the following description: a stream of numbers is random if, after observing N numbers being produced, you will not be able to arrive at a prediction of the $(N + 1)^{\text{th}}$ number to better than that given by its *probability* (for discrete random numbers), or to a prediction of its falling inside a certain interval better than given by its *probability density* (for continuous random numbers). That is, 'you cannot beat the house'.

¹²See, *e.g.* Signor Aglié's discussion of the dimensions of a newspaper stand in U. Eco, *Il Pendolo di Foucault*. Another example is the discovery of Dr. Irving Joshua Matrix that the decimals of π , correctly interpreted, contain the complete history of the human race (as reported in an interview by Martin Gardner).

⌘ All definitions of random streams follow this approach if you study them closely. It ultimately leads to the frequentist interpretation of probability.

1.2.3 What is probability? Relying on a random stream

The (to my mind) most operationally useful definition of ‘probability’ resides in the following description: given a stream of truly random numbers it may be possible, after observing N numbers being produced, to determine *probabilities*, that is, the fraction of numbers that attain a certain value (for discrete random numbers), or the fraction of numbers ending up inside a predetermined interval (for continuous random numbers). That is, ‘the house will not beat you’, an act of faith that can only be vindicated once $N = \infty$ has been reached and we are all dead.

⌘ This is the *frequentist* interpretation. Among other ones are the *propensity* interpretation, which is untenable without becoming frequentist, and the Bayesian, that rather describes a methodology for obtaining the probabilities.

1.2.4 A difference between physics and mathematics

The above notions of randomness and probability are circular and based on an operational picture of computational practice. The *mathematical* branch of probability theory, on the other hand (based on σ -algebra’s and measures), is rigorous, but while it describes what you can do with probabilities, it never asks the question of what probability *means*.

⌘ The meaning of meaning is not mathematics.

1.3 Miscellaneous probability items

1.3.1 Some notation in these notes

When sums or integrals are given without limits they are understood to run over all applicable real values (from minus to plus infinity).

Probability will always refer to a probability *density*, never to the (cumulative) probability distribution beloved by mathematicians. This is because the notion of density is defined in any dimension, and distribution is not.

The logical step function $\theta(A)$ has for its argument a statement. $\theta(A)$ equals 1 if A is true, and 0 if A is false. If a and b are integers, $\theta(a = b)$ is the Kronecker delta.

The ‘falling power’ $N^{\underline{k}}$ is defined as $N!/(N-k)! = N(N-1) \cdots (N-k+1)$.

1.3.2 Moments and characteristic function

Given a probability density $P(\mathbf{x})$ with support Γ we define the expectation value of a function $\varphi(\mathbf{x})$ by

$$\langle \varphi \rangle_P \equiv \int_{\Gamma} d\mathbf{x} P(\mathbf{x}) \varphi(\mathbf{x}) ; \quad (1)$$

this is exactly what probability means¹³. The subscript P is left out when no confusion can arise. For a *one-dimensional* density $P(x)$, we define the k^{th} moment as $\langle x^k \rangle$. Useful notions are

$$\begin{aligned} \text{the mean} & : \langle x \rangle , \\ \text{the variance} & : \sigma(x)^2 \equiv \langle x^2 \rangle - \langle x \rangle^2 , \\ \text{the characteristic function} & : \chi(z) = \chi_P(z) = \langle \exp(izx) \rangle . \end{aligned} \quad (2)$$

The last is defined if all moments are finite. $\sigma(x)$ is called the *standard deviation*. Obviously, $\chi(0) = 1$, $\chi'(0) = i \langle x \rangle$, $\chi''(0) = -\langle x^2 \rangle$, and

$$P(y) = \frac{1}{2\pi} \int dz \chi_P(z) e^{-izy} . \quad (3)$$

⋈ In the following we will generally assume that all moments exist.

1.3.3 The Chebyshev-Bienaymé inequality

Consider a (one-dimensional) probability density $P(x)$ with finite mean $\langle x \rangle = m$ and variance. Then, for any $a > 0$,

$$\begin{aligned} \sigma(x)^2 &= \int dx P(x) (x - m)^2 \geq \int_{|x-m|>a} dx P(x) (x - m)^2 \\ &\geq \int_{|x-m|>a} dx P(x) a^2 = a^2 \text{Prob}(|x - m| > a) . \end{aligned} \quad (4)$$

¹³In the frequentist sense.

We see that the probability for x to fall further from its mean than k times its standard deviation $\sigma(x)$ is always less than $1/k^2$. Since the estimated variance, E_2 , decreases as $1/N$ this theorem guarantees that Monte Carlo integration converges as long as the finiteness of $\langle E_2 \rangle$ is established with some confidence.

⌘ In the sense of this theorem, the probability density

$$P(x) \propto \theta(|x - m| \geq \sigma) \left(\frac{\sigma}{|x - m|} \right)^{2+\epsilon}$$

where ϵ is finite but as small as you like, is the widest possible density.

1.3.4 The Central Limit Theorem

Let the real numbers $x_{1,2,\dots,n}$ be iid random with density $P(x)$, and we assume that all the moments are finite. We denote $m = \langle x_j \rangle$, $\sigma = \sqrt{\sigma(x_j)^2}$. Then,

$$\xi = \frac{1}{n} \sum_{j=1}^n x_j \tag{5}$$

is also a random variate, with density $P_n(\xi)$. Its characteristic function is

$$\chi_{P_n}(z) = \langle \exp(iz\xi) \rangle = \prod_{j=1}^n \langle \exp(izx_j/n) \rangle = \chi_P(z/n)^n, \tag{6}$$

where $\chi_P(z)$ is the characteristic function of the x_j . We can approximate, for large n ,

$$\begin{aligned} \log(\chi_{P_n}(z)) &= n \log(\chi_P(z/n)) \\ &= n \log \left(1 + iz \langle x \rangle / n - z^2 \langle x^2 \rangle / 2n^2 + \mathcal{O}(1/n^3) \right) \\ &= izm - z^2 \sigma^2 / 2n + \mathcal{O}(1/n^2). \end{aligned} \tag{7}$$

Thus, for large n , we have approximately

$$\begin{aligned} P_n(\xi) &\approx \frac{1}{2\pi} \int dz \exp \left(iz(m - \xi) - \frac{z^2 \sigma^2}{2n} \right) \\ &= \sqrt{\frac{n}{2\pi\sigma^2}} \exp \left(\frac{-n(\xi - m)^2}{2\sigma^2} \right) \end{aligned} \tag{8}$$

This is the simplest version of the Central Limit Theorem: the distribution of the average of n iid¹⁴ random numbers with mean m and standard deviation σ approaches a Gaussian with mean m and standard deviation σ/\sqrt{n} .

⌘ This relies on the finiteness of the moments. A counterexample is the density $P(x) = (1 + x^2)^{-1}/\pi$ for which $P_n(\xi) = P(\xi)$ for *any* n . Note, however, that in that case even the mean is not well defined, and the variance is not finite.

¹⁴Independent, identically distributed.

2 Monte Carlo integration

2.1 The Monte Carlo idea

2.1.1 Point sets and expectation values

In standard Monte Carlo integration¹⁵ the random numbers \mathbf{x} are assumed to be iid with a probability density $P(\mathbf{x})$ that has support Γ . The integrand is $f(\mathbf{x})$, and the *weight* $w(\mathbf{x})$ is given by

$$w(\mathbf{x}) = f(\mathbf{x})/P(\mathbf{x}) \quad . \quad (9)$$

The weights w are also iid. The expectation values J_n are given by

$$J_n = \langle w \rangle_P = \int_{\Gamma} d\mathbf{x} P(\mathbf{x}) w(\mathbf{x})^n \quad , \quad (10)$$

and therefore $J_0 = 1$ and

$$J_1 = \int_{\Gamma} d\mathbf{x} f(\mathbf{x}) \quad , \quad (11)$$

which is the sought-after integral.

2.1.2 Integration as archetype

In a sense *any* Monte Carlo calculation, that is any calculation the outcome of which depends on at least one random number is an integration. Because if the outcome \mathcal{R} of a calculation depends on $N \geq 1$ random numbers, $\mathcal{R} = \mathcal{F}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, then its expected value is nothing but an integral,

$$\langle \mathcal{R} \rangle = \int d\mathbf{y}_1 \cdots d\mathbf{y}_N P(\mathbf{y}_1) \cdots P(\mathbf{y}_N) \mathcal{F}(\mathbf{y}_1, \dots, \mathbf{y}_N) \quad . \quad (12)$$

It is therefore sensible to concentrate on Monte Carlo integration.

⌘ This is strictly formal; any serious Monte Carlo simulation easily employs many millions of random numbers, and its numerical result is therefore a multimillion-dimensional integral. Nevertheless the above point of view is useful to keep in mind.

¹⁵As opposed to Quasi-Monte Carlo.

2.1.3 Point sets, ensembles, and the Leap Of Faith

In the Monte Carlo approach we employ a *point set* \mathbf{X} consisting of N points \mathbf{x}_j : $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ where the \mathbf{x} 's are sampled from the distribution $P(\mathbf{x})$. An individual point \mathbf{x}_j is called an *event*, and $w_j \equiv w(\mathbf{x}_j)$ is called the corresponding *event weight*. Given \mathbf{X} we can compute

$$S_k = S_k(\mathbf{X}) = \sum_{j=1}^N w_j = \sum_{j=1}^N w(\mathbf{x}_j)^k \quad (k = 0, 1, 2, 3, 4) \quad . \quad (13)$$

Given the notion of a stream of random numbers, we envisage a great number ($\rightarrow \infty$) of point sets \mathbf{X} : the *ensemble of point sets*. Averaging over the random numbers is averaging over the ensemble, which is defined by the combined probability density it imposes on the points. For instance, for point sets defined on the d -dimensional hypercube $I^d = (0, 1)^d$ it is simply

$$P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = 1 \quad . \quad (14)$$

This immediately shows the uniformity and the iid property. An important thing to remember is that in any calculation we *assume* that \mathbf{X} is a ‘typical’ member of the ensemble, and that therefore the ensemble averages are meaningful for the given point set. This is the Leap of Faith.

⌘ The Leap of Faith can, and if possible should, be vindicated by repeating a computation with a different point set, obtained by either taking a different part of the generated random number stream, or switching to another random number generator.

2.2 Finding estimators

2.2.1 Direct sums, unequal-sums, and expectation values

We consider N iid random numbers w_j ($j = 1, 2, \dots, N$) with expectation values

$$\langle w_j^k \rangle \equiv J_k \quad . \quad (15)$$

We define the direct sum S_m ($m = 0, 1, 2, \dots$) as

$$S_m = \sum_{j=1}^N (w_j)^m \quad : \quad (16)$$

these sums are computable in *linear* time, $\mathcal{O}(N)$. S_0 is simply equal to N . The unequal-sums S_{m_1, m_2, \dots, m_k} are defined as

$$S_{m_1, m_2, \dots, m_k} = \sum_{j_1, j_2, \dots, j_k=1}^N (w_{j_1})^{m_1} (w_{j_2})^{m_2} \dots (w_{j_k})^{m_k} \quad (17)$$

with the constraint that the indices j_1, \dots, j_k are all *different* from one another. The unequal-sum S_{m_1, m_2, \dots, m_k} therefore contains $N^{\underline{k}}$ terms. Its straightforward computation takes time $\mathcal{O}(N^k)$. By the iid assumption we have

$$\langle S_m \rangle = N J_m \quad , \quad \langle S_{m_1, m_2, \dots, m_k} \rangle = N^{\underline{k}} J_{m_1} J_{m_2} \dots J_{m_k} \quad . \quad (18)$$

The falling powers $N^{\underline{k}}$ are discussed in appendix 11.0.1. We can relate products of direct sums to combinations of unequal-sums by the following rule:

$$\begin{aligned} S_{m_1, m_2, \dots, m_k} S_p &= S_{m_1+p, m_2, \dots, m_k} + S_{m_1, m_2+p, \dots, m_k} + \dots \\ &\dots + S_{m_1, m_2, \dots, m_k+p} + S_{m_1, m_2, \dots, m_k, p} \quad . \end{aligned} \quad (19)$$

Explicit relations are given in appendix 11.0.2.

⌘ This provides a way to evaluate unequal-sums in linear time, and a way to find expectation values of nonlinear combinations of sums. Conversely, we can find the combination of direct sums that has a given expectation value.

2.2.2 The Monte Carlo estimators

Using the direct sums S_k we define three *Monte Carlo estimators* :

$$\begin{aligned} E_1 &= \frac{1}{N} S_1 \quad , \\ E_2 &= \frac{1}{N N^{\underline{2}}} (N S_2 - S_1^2) \quad , \\ E_4 &= \frac{N^{\underline{2}}}{N^3 N^{\underline{4}}} (N S_4 - 4 S_3 S_1 + 3 S_2^2) \\ &\quad + \frac{1}{N^4} \left(\frac{2}{N^2 N^{\underline{2}}} - \frac{4}{N^3} \right) (N S_2 - S_1^2)^2 \quad . \end{aligned} \quad (20)$$

Since the point set \mathbf{X} is an element of the ensemble of point sets, the numbers $E_{1,2,4}$ are also (one-dimensional) random numbers with their own mean and

variance¹⁶. Using 11.0.2 we can prove

$$\begin{aligned}
\langle E_1 \rangle &= J_1 , \\
\langle E_2 \rangle &= \sigma(E_1)^2 = \frac{1}{N} (J_2 - J_1^2) , \\
\langle E_4 \rangle &= \sigma(E_2)^2 = \frac{1}{N^3} (J_4 - 4J_3J_1 + 3J_2^2) \\
&\quad + \left(\frac{2}{N^2N^2} - \frac{4}{N^3} \right) (J_2 - J_1^2)^2 . \tag{21}
\end{aligned}$$

We see that (in an ensemble sense) $\langle E_1 \rangle$ is the desired integral J_1 . E_2 informs about the ensemble variance in the probability density of E_1 values, and it has its own ensemble variance, estimated by E_4 . Note that $\langle E_2 \rangle < \infty$ if the integrand is *quadratically integrable*, but $\langle E_4 \rangle < \infty$ only if the integrand is *quartically integrable*.

⌘ There is also an E_8 with $\langle E_8 \rangle = \sigma(E_4)^2$, and so on. These ever more complicated estimators are not very relevant. E_1 and E_2 are well known, and E_4 *deserves to be*.

2.2.3 Positivity of E_2 and E_4

In the Monte Carlo integral we have $\langle w(\mathbf{x}) \rangle = J_1$. Now write $u(\mathbf{x}) = w(\mathbf{x}) - J_1$. Then we have

$$\begin{aligned}
J_2 - J_1^2 &= \int d\mathbf{x} P(\mathbf{x}) u(\mathbf{x})^2 , \\
J_4 - 4J_3J_1 + 3J_2^2 - 4(J_2 - J_1^2)^2 &= \\
\frac{1}{2} \int d\mathbf{x} d\mathbf{y} P(\mathbf{x}) P(\mathbf{y}) (u(\mathbf{x})^2 - u(\mathbf{y})^2)^2 & \tag{22}
\end{aligned}$$

so both E_2 and E_4 have positive expectation value¹⁷. For a given \mathbf{X} , define $u_j = w(\mathbf{x}_j) - J_1$. Then

$$NS_2 - S_1^2 = \frac{1}{2} \sum_{j,k=1}^N (u_j - u_k)^2 \geq 0 . \tag{23}$$

¹⁶In the sense that every new chosen point set \mathbf{X} yields its own values for $E_{1,2,4}$.

¹⁷As they should!

So the estimator E_2 will always result in a positive number¹⁸. However, suppose that w_j only takes the values 0 and 1, so that $S_k = Nb$ for all $k > 0$, where $0 < b < 1$. Then, E_4 evaluates to a negative number if $b(1 - b) > 1/4 - (N - 2)/2N(4N - 6)$, so E_4 cannot be guaranteed to be nonnegative. On the other hand,

$$N^2 (NS_4 - 4S_3S_1 + 3S_2^2) - 4 (NS_2 - S_1^2)^2 = \frac{N^2}{2} \sum_{j,k=1}^N (u_j^2 - u_k^2)^2 \geq 0 . \quad (24)$$

This suggests a slight modification of the estimators $E_{2,4}$.

2.3 Estimating in practice

2.3.1 Improved estimators

We can use the following ‘improved’ estimators:

$$\begin{aligned} E_1 &= \frac{1}{N} S_1 , \\ \hat{E}_2 &= \frac{1}{N^3} (NS_2 - S_1^2) , \\ \hat{E}_4 &= \frac{1}{N^7} \left(N^2 (NS_4 - 4S_3S_1 + 3S_2^2) - 4 (NS_2 - S_1^2)^2 \right) . \end{aligned} \quad (25)$$

These are equal to the exact ones up to $1/N$ corrections, and are nonnegative by construction.

⋈ Since in any serious calculation N is of the order of at least a few thousand, a $1/N$ correction in uncertainty estimates is not a big deal.

2.3.2 Numerical stability; the extended CGL algorithm

Computing $E_{2,4}$ involves large cancellations. Even E_2 can almost never be computed using single precision if Eq.(25) is used¹⁹. The following algorithm avoids this. Suppose $n - 1$ random numbers have been used already, and the contribution w_n of the n^{th} is applied. Define

$$U_k(n) = \frac{1}{n} \sum_{j=1}^n w_j^k \quad (26)$$

¹⁸Barring numerical accidents.

¹⁹This is really true, as *anyone* who tried it seriously will testify.

and

$$\begin{aligned}
M(n) &= U_1(n) , \\
P(n) &= U_2(n) - U_1(n)^2 , \\
Q(n) &= U_3(n) - 3U_2(n)U_1(n) + 2U_1(n)^3 , \\
R(n) &= U_4(n) - 4U_3(n)U_1(n) + 3U_2(n)^2 - 4P(n)^2 ,
\end{aligned} \tag{27}$$

and also

$$m = M(n-1) , \quad p = P(n-1) , \quad q = Q(n-1) , \quad r = R(n-1) . \tag{28}$$

Let us also define $u = w_n - m$. Then we can update in linear time as follows:

$$\begin{aligned}
M(n) &= m + u/n , \\
P(n) &= \frac{n-1}{n} \left(p + \frac{u^2}{n} \right) , \\
Q(n) &= \frac{n-1}{n} \left(q + \frac{n-2}{n^2} u^3 - \frac{3p}{n} u \right) , \\
R(n) &= \frac{n-1}{n} \left(r + \frac{1}{n} \left(p - \frac{n-2}{n} u^2 \right)^2 - 4 \left(\frac{q}{n} u - \frac{p}{n^2} u^2 \right) \right) .
\end{aligned} \tag{29}$$

This algorithm is useful since (a) it can be used to update in constant time, and (b) it is free of large cancellations.

⋈ The use of M and P in the computation of \hat{E}_2 is the original CGL algorithm [5]. The Q and R necessary for \hat{E}_4 is discussed in Bakx *et al.* [6].

2.3.3 How to *do* Monte Carlo integration

The recommended way to compute a Monte Carlo estimate of the integral

$$J_1 = \int_{\Gamma} d\mathbf{x} f(\mathbf{x}) \tag{30}$$

is then as follows: generate a point set \mathbf{X} of iid random numbers \mathbf{x}_j ($j = 1, \dots, N$) with density $P(\mathbf{x})$. Every time a new point \mathbf{x}_n is added, compute $w(\mathbf{x}_n) = f(\mathbf{x}_n)/P(\mathbf{x}_n)$. Then, update M, P, Q , and R . After each point, the running (‘improved’) estimates are given by

$$E_1 = M(n) , \quad \hat{E}_2 = P(n)/n , \quad \hat{E}_4 = R(n)/n^3 . \tag{31}$$

Below we give the algorithm in pseudocode

Algorithm 1 Extended CGL algorithm for numerically safe updating of running Monte Carlo estimators

{ At every call to this algorithm an event weight w is inputted. The number n is the number of event weights inputted so far. The estimators $\hat{E}_{1,2,4}$ are the output. The numbers M, P, Q, R are kept internally. }

if $n = 0$ **then**
 $[M, P, Q, R] \leftarrow [0, 0, 0, 0]$ {Initialization}
end if
 $[m, p, q, r] \leftarrow [M, P, Q, R]$ {Variables used in the update}
 $n \leftarrow n + 1$
 $u \leftarrow w - m$
 $M \leftarrow m + u/n$
 $P \leftarrow (n - 1)(p + u^2/n)/n$
 $Q \leftarrow (n - 1)(q + (n - 2)u^3/n^2 - 3pu/n)/n$
 $R \leftarrow (n - 1)(r + (p - (n - 2)u^2/n)^2/n - 4(qu/n - pu^2/n^2))/n$
 $[\hat{E}_1, \hat{E}_2, \hat{E}_4] \leftarrow [M, P/n, R/n^3]$ {The estimators so far}

⋈ The updating allows for monitoring how the calculation progresses. This is a very important technique to gauge the quality of the computation, as we shall see.

2.3.4 How to *report* Monte Carlo integration

After a point set \mathbf{X} of N points has been used to do a Monte Carlo integral, the estimate of the integral is E_1 . The estimate of its ensemble variance is \hat{E}_2 , and the estimate of the ensemble variance of E_2 is \hat{E}_4 . The best way to report the result of the computation is

$$J_1 = E_1 \pm \left(E_2^{1/2} \pm E_4^{1/4} \right) . \quad (32)$$

2.3.5 How to *interpret* Monte Carlo integration

The ‘Monte Carlo error’ is given by $\hat{E}_2^{1/2}$. If the Central Limit Theorem holds, this gives the Gaussian *confidence levels* associated with the answer. The ‘error on the error’ $\hat{E}_4^{1/4}$ informs about the uncertainty in the confidence levels, which can be quite important.

We have, for a given integrand,

$$\frac{\hat{E}_2^{1/2}}{E_1} \sim \frac{1}{N^{1/2}} \quad . \quad (33)$$

This is the well-known, slow but universal²⁰ behaviour of Monte Carlo integration. Not so well known but important is the relative uncertainty of the error:

$$\frac{\hat{E}_4^{1/4}}{\hat{E}_2^{1/2}} \sim \frac{1}{N^{1/4}} \quad . \quad (34)$$

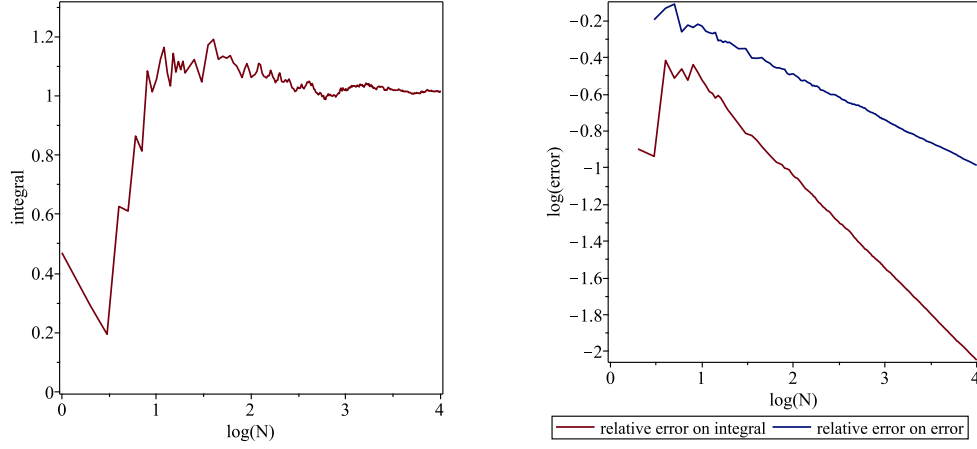
The convergence of the error estimate is much slower than that of the integral. The error can be translated (in the Central Limit Theorem sense) into Gaussian confidence intervals. Here the error on the error becomes important. The 1σ confidence level (two-sided) is 68.2%; if the relative error on the error is 30 per cent (admittedly quite large) the confidence level can range from 51.6% to 80.6%. It is important to compute \hat{E}_4 even if only to confirm that the error on the error is small.

⌘ The error on the error is not widely known, unfortunately. Experience shows that it can be uncomfortably large.

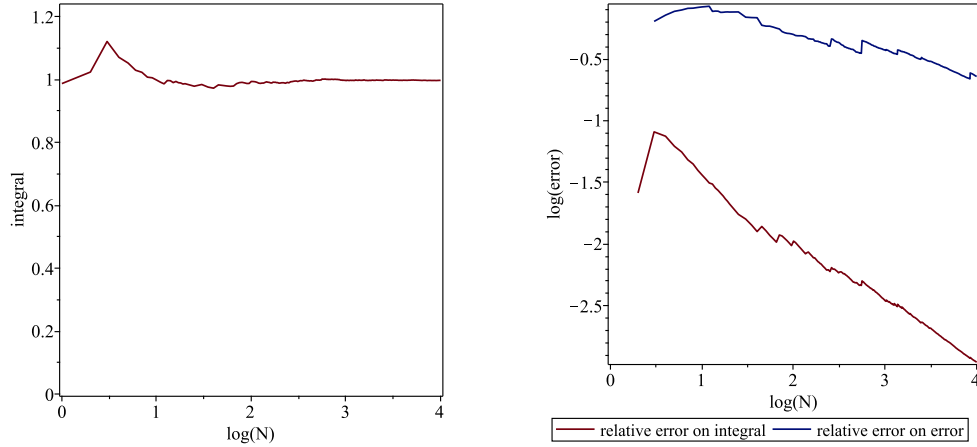
2.4 A test case

We consider the MC integration of the function $f(a; x) = (1 + a)x^a\theta(0 < a \leq 1)$, with $a > -1$, for various decreasing values of a . The *real* value of the integral is unity. We use the same set of 10000 (pseudo)random numbers for all plots.

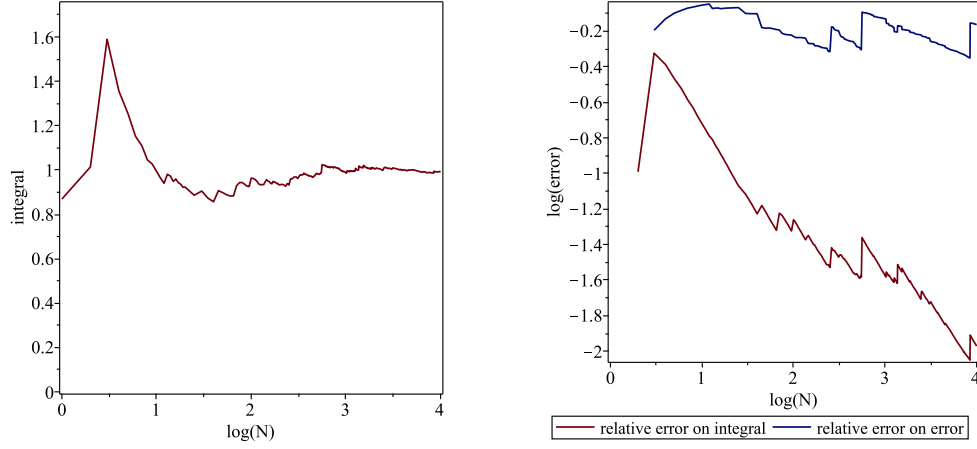
²⁰In particular, independent of the dimensionality of Γ .



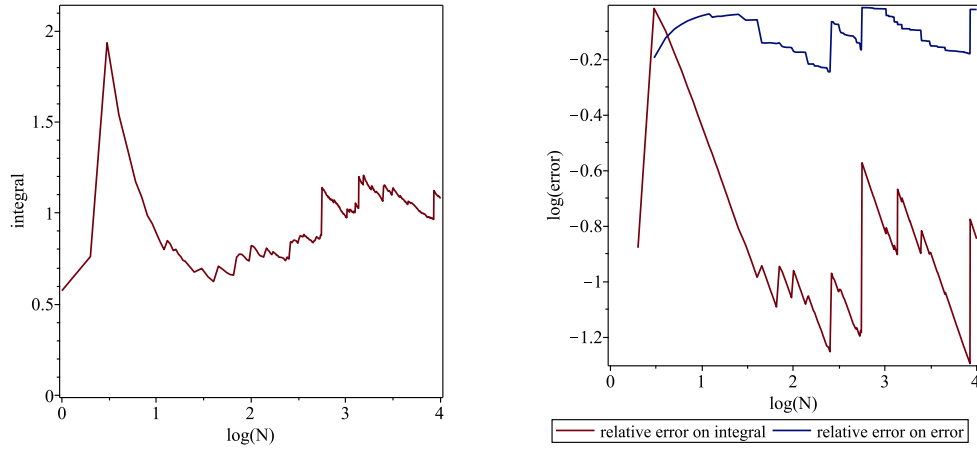
For $a = 2$ the integrand is bounded and perfectly integrable. The estimates of the relative first- and second-order error decrease as $N^{-1/2}$ and $N^{-1/4}$ respectively.



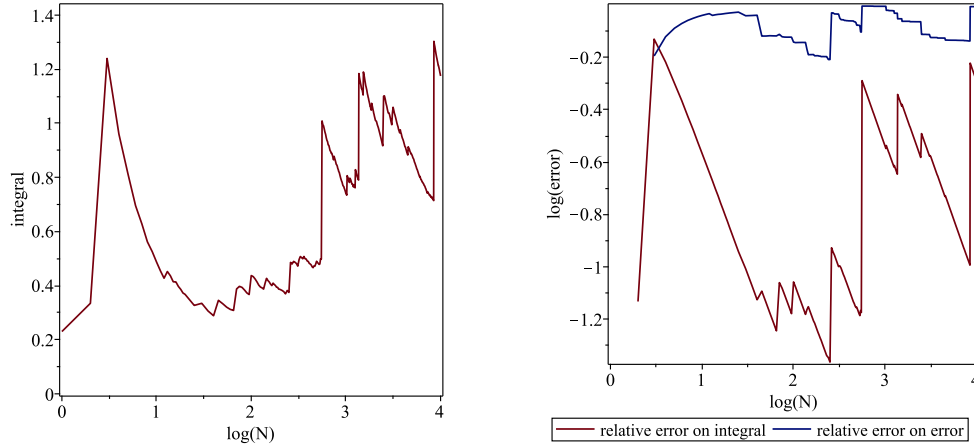
For $a = -0.1$ the integrand is no longer bounded but integrable and also square and quartically integrable. The first-order error still decreases smoothly, but the second-order error remains quite large (22% at $N = 10000$). The small jumps in the curves arise from MC points close to zero.



For $a = -0.4$ the integrand is no longer quartically integrable, witness the sizeable jumps in the error curves. The second-order error cannot really be said to decrease at all even though the estimated error at $N = 10000$ is still a reasonable 1%, but its own error is around 60% (relative).



For $a = -0.7$ the integrand is no longer square integrable, and the first-order error estimate cannot be trusted at all, from the fact that here $E_4^{1/2}/E_2$ is around unity.



For $a = -0.9$ integrability is almost completely lost: E_2 jumps all over the place and even E_1 cannot be said to converge

From these simple tests a number of conclusions can be drawn. In the first place, integration problems like these show, in the development of the estimators, the phenomenon of ‘*punctuated equilibrium*’: the *a priori* behaviour of E_2 and E_4 is still as $1/N$ and $1/N^3$, respectively, but interspersed with jumps that become more pronounced as the integrand becomes more and more singular²¹. Eventually these jumps counteract the decreasing of the estimators. In the second place, we see that in any MC integration one should not be content with the final numbers alone, but the integration should be monitored while it proceeds. Jumps signal singularities!

2.5 Exercises

Exercise 1 The CGL algorithm

Write your own Monte Carlo evaluator. This should be a code in which (after initialisation) weights can be inputted one after the other, while a continuous estimate of $\hat{E}_{1,2,4}$ is updated using the CGL algorithm. That is, after a number of weights have been inputted, the code should be able to give the Monte Carlo estimators at any moment during the computation.

²¹A similar notion from a quite distant field is found in [7].

Exercise 2 Studying the test case

Using the random number generator of your choice, perform your own study of the test case of sect.2.4 by using the code of the previous exercise, and playing around with the parameter a . You may also want to consider some other function with an adjustable amount of singularity.

Exercise 3 For diehard combinatoricists only

Try to find the variance of \hat{E}_4 , and find the estimator E_8 that gives ‘the uncertainty in the error estimate on the error’. It involves things like $S_{1,1,1,1,1,1,1,1}$ and $S_{3,2,2,1} \dots$

3 Random number generation

3.1 Introduction to random number sources

3.1.1 Natural *vs* Pseudo

Random streams can be obtained from natural processes, or by computer. In the last case there *is* an algorithm and the ‘next’ number is entirely predictable²². Such streams aim not at *being* random but at *appearing* (sufficiently) random, and are called *pseudorandom* numbers streams²³. Both methods have advantages and drawbacks.

1. Natural random number streams: electronic noise, chaotic lasers, air pressure fluctuations, photonics, . . .
 - Advantages: unpredictability
 - Drawbacks: unrepeatability (in all cases), speed (in some cases), unguaranteed probability density (in all cases)

The picture shows IDQuantique’s Quantis-USB-4M Quantum Random Number Generator²⁴. It produces ‘truly random’ bits by observing indi-



vidual photons that are or are not reflected by a half-reflecting mirror. I timed (by hand) approximately 21 seconds to generate 80 Mbit and about the same time for 2×10^4 10-digit integers or 10-digit floating point numbers in $(0, 1)$. For such scientific tasks as Monte Carlo integration this is typically much too slow.

Its main application is cryptographical; there, true unpredictability and unrepeatability are essential.

²²By those who are in on the secret, or on the code.

²³Greek $\psi\epsilon\upsilon\delta\epsilon\iota\nu$, ‘to lie’.

²⁴<https://www.idquantique.com>

⌘ *I think that the direct use of a physical supply of random digits is absolutely unacceptable for this reason and for this reason alone.*

John von Neumann on the irreproducibility of natural random numbers [25].

2. Pseudorandom number streams: computer algorithms

- Advantages: understandability, repeatability
- Drawbacks: predictability, speed (in some cases), nonuniformity (in some cases)

⌘ *Anyone who considers arithmetical methods of producing random digits is, of course, in a state of sin.*

John von Neumann on the predictability of pseudo-random numbers [25].

3.2 Pseudo-random number streams

3.2.1 The structure of pseudo-random number algorithms

Consider algorithms working on the set of integers $a_j \in \{1, 2, \dots, M\}$, $j = 1, 2, 3, \dots$. Pseudorandom number generators²⁵ always generate the next number as a deterministic function of one or more previous ones. We shall first examine the algorithm

$$a_{n+1} = f(a_n) \quad . \quad (35)$$

Given a starting value a_1 it generates the stream

$$a_2 = f(a_1) \quad , \quad a_3 = f(a_2) \quad , \quad f(a_4) = f(a_3) \quad , \dots \quad (36)$$

If $f(a_L) = a_p \in \{a_1, a_2, a_3, \dots, a_L\}$ (*i.e.* a number comes up that was already generated before) the series $(a_p, a_{p+1}, \dots, a_L)$ will start repeating itself and pseudorandomness is lost. L is the *lifetime* of the algorithm f for starting value a_1 . The series $(a_1, a_2, \dots, a_{p-1})$ can be called the *runup*, and the series (a_p, \dots, a_L) the *cycle*.

²⁵Commonly denoted as PRNGs.

3.2.2 The set of all algorithms

The above algorithm f can be *completely* specified by the list

$$S_f = \left[f(1), f(2), f(3), \dots, f(M-1), f(M) \right] . \quad (37)$$

Many different-looking algorithms can end up with the same S_f . A few conclusions immediately follow:

1. There are precisely M^M possible algorithms of the form (35).
2. The maximum possible lifetime is M .
3. The number of algorithms with the maximum lifetime for a *given* starting value is $M!$.
4. The number of algorithms with the maximum lifetime for *any* starting value is $(M-1)!$, the number of permutations with no subcycles. These have no runup.

3.2.3 The concept of an arbitrary algorithm

We can pick a ‘arbitrarily chosen algorithm’ by choosing S_f randomly out of the M^M possibilities, and so arrive at probabilistic statements about ‘arbitrary’ algorithms²⁶. The probability to have $a_2 = f(a_1) \neq a_1$ is $1 - 1/M$; the probability to have $a_3 = f(a_2) \neq a_{1,2}$ is $1 - 2/M$, and so on: the probability of having a lifetime L of *at least* k is therefore²⁷

$$\text{Prob}(L \geq k) = \frac{M^k}{M^k} . \quad (38)$$

The probability of having a lifetime of precisely k is

$$\text{Prob}(L = k) = \text{Prob}(L \geq k) - \text{Prob}(L \geq k+1) = \frac{k M^k}{M^{k+1}} . \quad (39)$$

The expected lifetime given a fixed arbitrary starting value a_1 is

$$\sum_{k=1}^M k \text{Prob}(L = k) = \sum_{k \geq 1} k \text{Prob}(L \geq k) - \sum_{k \geq 1} k \text{Prob}(L \geq k+1)$$

²⁶The phrasing ‘random’ algorithm may lead to confusion.

²⁷Correctly, this gives a probability zero for a lifetime longer than M .

$$\begin{aligned}
&= \sum_{k \geq 1} \text{Prob}(L \geq k) = \sum_{k \geq 1} \frac{M^k}{M^k} \\
&= \int_0^\infty dx e^{-x} \sum_{k \geq 1} \frac{x^k M^k}{k! M^k} \\
&= \int_0^\infty dx e^{-x} \left(-1 + \left(1 + \frac{x}{M} \right)^M \right) \\
&= \sqrt{\frac{\pi M}{2}} - \frac{1}{3} + \sqrt{\frac{\pi}{288 M}} + \mathcal{O}\left(\frac{1}{M}\right) , \tag{40}
\end{aligned}$$

where the last line is derived in appendix [11.0.4](#).

3.2.4 Lessons from random algorithms

A number of probabilistic remarks hold for randomly chosen algorithms.

- The number M should be very large, since true random number streams correspond to $M \rightarrow \infty$.
- The probability of picking an algorithm with lifetime $L = M$ is $M!/M^M \approx \sqrt{2\pi M} \exp(-M)$, *i.e.* extremely small.
- The expected lifetime is of order \sqrt{M} , so very much smaller than the maximum lifetime.
- The probability of occurrence of at least one *selfie*, a number a such that $f(a) = a$, is *not* small, around 63% (see appendix [11.0.5](#)).

We can draw a number of inferences that are important for practical PRNGs.

- One should aim for long lifetimes. This is of course only the very crudest requirement on a PRNG.
- Good algorithms are rare, and ‘lie close to’ bad or mediocre algorithms.
- Rounding errors are to be avoided²⁸. This is why in many cases the internal arithmetic is done with integer a_j , and only at the output stage the floating-point number $a_j/M \in (0, 1]$ is returned.

²⁸Rounding errors mean that the algorithm does not do *exactly* what it is meant to do, and so deviates from a possible good algorithm into the jungle of bad algorithms.

- For an algorithm with maximum lifetime, at the end of the lifetime *all* numbers will have been generated exactly *once*. This does not look random at all. One should avoid coming close to exhausting the whole lifetime. In practice, if you envision using n random numbers, then L should be at least of order n^2 .

3.2.5 Shift-register PRNGs

One can let the next number a_n depend not only on a_{n-1} but on a register of size r :

$$a_n = f(R_n) \quad , \quad R_n = (a_{n-1}, a_{n-2}, \dots, a_{n-r+1}, a_{n-r}) \quad . \quad (41)$$

The *next* number will be

$$a_{n+1} = f(R_{n+1}) \quad , \quad R_{n+1} = (a_n, a_{n-1}, \dots, a_{n-r+2}, a_{n-r+1}) \quad ; \quad (42)$$

a_n has entered from the left, and a_{n-r} has dropped out on the right : hence the name *shift-register* PRNG. This is in fact just an artifice to enlarge M since we can map the register to integers by

$$R_n \leftrightarrow 1 + (a_{n-1} - 1) + M(a_{n-2} - 1) + M^2(a_{n-3} - 1) + \dots + M^{r-1}(a_{n-r} - 1) \quad , \quad (43)$$

which is an integer between 1 and M^r . The use of the register is just a way of writing very large numbers in base M , and only using the leading ‘digit’ as the random number. The maximum possible lifetime is now M^r .

⌘ Shift-register PRNGs can be deceptively simple : an algorithm like $a_n \sim a_{n-s} \pm a_{n-r}$ ($1 \leq s < r$) is often already quite good.

3.3 Bad and good algorithms

3.3.1 Bad: the midsquare algorithm

This is one of the oldest attempts at a PRNG algorithm. Of the square of a (not too small) integer, the *last* digit(s) are easily predictable, as are the first few digits. By taking squares of $2d$ -digits numbers, and retaining the middle $2d$ digits of the result (adding zeroes if leading digits in the $4d$ -digit square

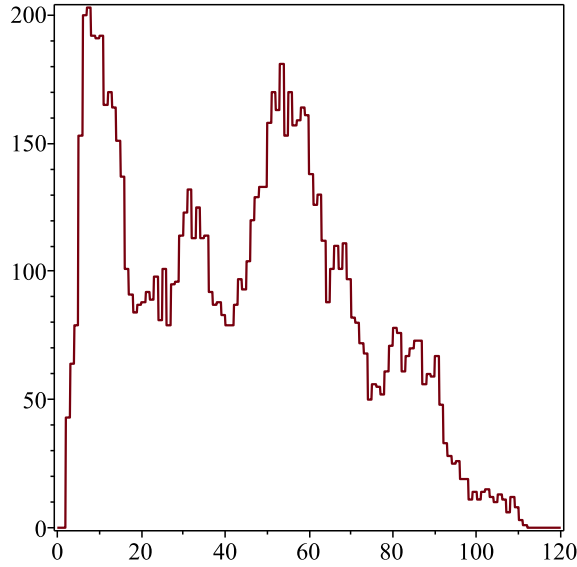
are missing), one hopes to achieve a ‘difficultly-predictable’ sequence. If K is the maximum integer as before it would read, using the ‘floor’ function,

$$a_n = \left\lfloor a_{n-1}^2 / \sqrt{K} \right\rfloor \bmod K . \quad (44)$$

Typically K would have to be a perfect square²⁹. For $K = 100$ this would give, for instance the sequence

$$63, 96, 21, 44, 93, 64, 9, 8, 6, 3, 0 .$$

Direct inspection tells us that, for $K = 100$: (i) there are 4 selfies, to wit 00, 10, 50, and 60 ; (ii) 61 sequences end in 0 ; (iii) the largest lifetime is 15, achieved for the starting value 42 (!) ; (iv) the average lifetime is 5.76, way below the ‘expected’ value of about 12.



Here we give the distribution of lifetimes for $K = 10^4$ for all starting values from 1 to 9999. The maximum lifetime is 111 (achieved for the starting value 6239), not even equal to the expected value of around 125. The mid-square method can be considered a typical example of an ‘arbitrary’ algorithm.

⌘ The midsquare method with $K = 10^{10}$ has been used with broadly satisfactory results in the early 1950’s^a but should be considered totally obsolete for modern applications.

^aAs reported by P.C.Hammer in [4], p33. A lifetime of at least 10^4 was found for starting value 111111111.

²⁹This is not strictly necessary, since the floor function erases rounding errors in the \sqrt{K} unless K is really huge.

3.3.2 Bad: the chaos approach and the logistic map

The chaotic behaviour of certain dynamical systems would seem to be a source of ‘random’, or at least unpredictable, sequences. A good example is provided by the fully chaotic *logistic map*. This is based on real numbers rather than integers:

$$x_n = 4 x_{n-1} (1 - x_{n-1}) \quad . \quad (45)$$

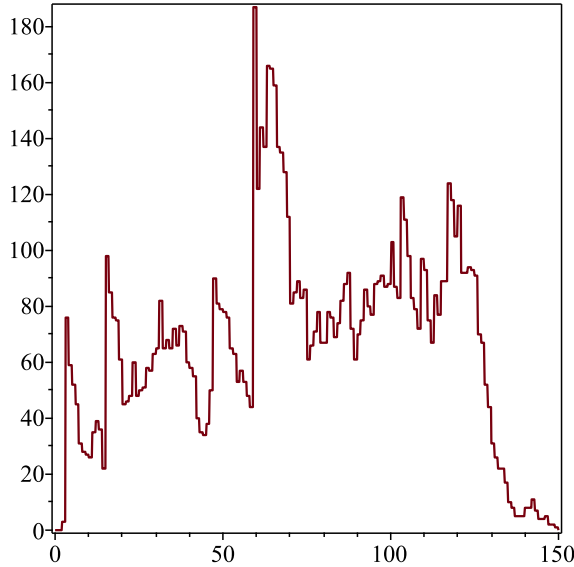
It is easy to see that this will almost certainly yield a sequence with infinite lifetime. If we write

$$x_n = \sin^2 \left(\frac{\pi}{2} y_n \right) \quad , \quad 0 \leq y_n \leq 1 \quad , \quad (46)$$

The map (45) corresponds to

$$y_n = 2y_{n-1} \mod 1 \quad . \quad (47)$$

Now, every non-rational number in $(0, 1)$ has a binary expansional that does not repeat. Therefore, the logistic map (45) will have infinite lifetime — in mathematical principle. However, we can only use finite-precision floating-point numbers. Therefore if y_n is given with (say) 100 binary digits’ precision, then iterating (47) would give $y_{n+101} = 0$, a selfie. Of course we do not use the y ’s but the x ’s but this means that the selfie $x = 0$ is only avoided because of rounding errors ! Consequently we expect the logistic map (in finite precision) to be an ‘arbitrary’ algorithm. The same holds for all dynamical systems that are chaotic : the very fact that their behaviour depends extremely sensitively on the initial conditions guarantees that in practice it is driven by rounding errors.



The lifetimes of the logistic-map algorithm for numbers restricted to 4 decimal digits, for all possible starting values. The maximum observed lifetime is 149, a little better than the expectation value 125.

⌘ An even stronger statement about binary (or decimal) expansions is possible. Almost all real numbers (in the sense of Lebesgue measures) are *normal*, that is their binary expansion contains every given block of digits (such as 0, 101, or 1001010110011) with asymptotically the correct frequency. This means that iterating (47) would give a perfect stream of random bits. Unfortunately nobody knows how to construct a normal number. π may be normal but that has not been proven yet.

3.3.3 Maybe not so bad: linear congruential methods

A very popular and simple (hence analysable!) algorithm is the *linear congruential* method. In terms of integers x_n ($n = 0, 1, 2, \dots$) it reads

$$x_{n+1} = (a x_n + c) \bmod m \quad (48)$$

The *modulus* m , the *multiplier* a , and the *increment* c are integers. One of the advantages is that the period (and other properties) can be determined [8]. The maximum lifetime (in this case, period) is obtained under the following conditions [9]:

1. c and m are relatively prime;
2. every prime factor of m divides $a - 1$;
3. if m is a multiple of 4, then so must $a - 1$ be.

When $c = 0$ the maximum period cannot be achieved³⁰. In that case the maximum period can be $m/4$ if m is a power of 2, or $m - 1$ if m is a prime, everything depending on the optimal choice of a [10]³¹.

It is interesting to note the following. Suppose $m = 2^n$, and $a = 2^p + q$, with $p > n/2$. Then

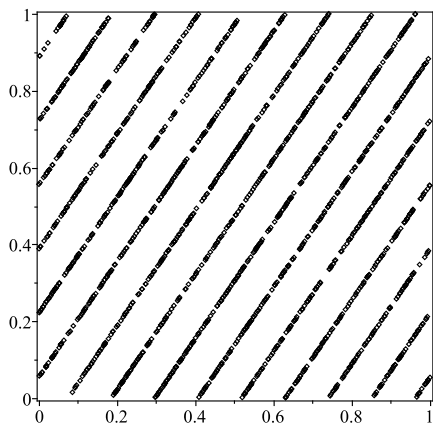
$$x_n = (2q x_{n-1} - q^2 x_{n-2}) \bmod m, \quad (49)$$

so this algorithm can be viewed as a shift-register PRNG as well. It also means that the triples $[x_{3k-2}, x_{3k-1}, x_{3k}]$ all lie on a 2-dimensional plane. Using the $\bmod m$ this plane intersects the 3-dimensional $m \times m \times m$ cube a number of times, so forming a collection of planes on which all generated triples must fall. This is easily extended to larger multiplets of points in hypercubes. Such structures of multiples of points are common to all multiplicative congruential PRNGs [11], the *maximum* number of planes in a k -dimensional hypercube being $m^{1/k}$. This limits the usefulness of multiplicative congruential generators, even when the period length m is acceptable.

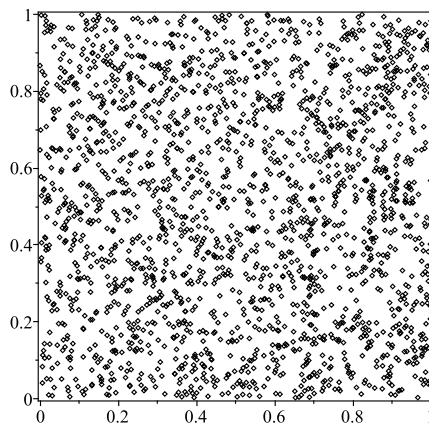
³⁰Since then $x_n = 0$ is a selfie.

³¹The indispensable reference here is [8].

3.3.4 Horrible: RANDU, and a possible redemption



Cube slice for RANDU



Cube slice for $a = 69069$.

For many years the PRNG called `RANDU` was popular[?]. It has $m = 2^{31}$, $c = 0$ and $a = 2^{16} + 3 = 65539$. Typically, x_1 is chosen to be 1. These values are mainly inspired by simplicity of implementation on a 32-bit machine. The choice of a has turned out to be about the worst possible; the number of planes (as discussed above) is very small in low dimensions. We can visualise this by, say, taking a slice of the cube with $0.32 \leq x_{3k} \leq 0.34$. `RANDU` gives only 15 planes! Notice that for such features to become obvious it was necessary for easy 2-d plotting to become available. A much better choice appears to be $a = 69069$ [13]³².

3.3.5 Good: RCARRY and RANLUX

An example of a shift-register PRNG that does better than a multiplicative congruential one, with a larger register, is `RCARRY`[14]. Its parameters are the modulus B , and two index parameters s and r , the latter being the register length. We use $B = 2^{24}$, $s = 10$ and $r = 24$. In addition there is a so-called *carry bit* c_n associated with the pseudorandom integer x_n . The complete register therefore reads

$$[x_{n-1}, \dots, x_{n-s}, \dots, x_{n-r}; c_{n-1}]$$

³² Viewed as shift-register PRNGs we have for `RANDU` $x_n = 6x_{n-1} - 9x_{n-2}$. For the redemption value $a = 69069$ we have $x_n = 7066x_{n-1} - 12482089x_{n-2}$. The coefficients are far larger, but the register length is still only 2.

and the algorithm is

$$x_n = (x_{n-s} - x_{n-r} - c_{n-1}) \bmod B \quad (50)$$

where $c_n = 0$ or $c_n = 1$ according to whether or not the $\bmod B$ was necessary. In order to avoid having to build the register anew every time, it is best to view it as cyclic with the starting point at the back.

Algorithm 2 The RCARRY algorithm using a cyclic register.

{The register is $[a_1, a_2, a_3, \dots, a_r]$, and its *last entry* is currently a_j . The current value of the carry bit is c . The new pseudorandom number is x/B .}

$p \leftarrow s + j$, if $p > r$ then replace $p \leftarrow p - r$ {cyclicity}

$q \leftarrow r + j$, if $q > r$ then replace $q \leftarrow q - r$ {cyclicity}

$y \leftarrow a_p - a_q - c$

if $y \geq 0$ **then**

$x \leftarrow y$ and $c \leftarrow 0$

else

$x \leftarrow x + B$ and $c \leftarrow 1$

end if

replace a_j by x

$j \leftarrow j - 1$; if $j = 0$ then replace $j \leftarrow r$ {cyclicity}

return x/B

As usual, the register has to be initialised. Best practice is to fill it with large ‘arbitrary’ numbers, of $\mathcal{O}(B)$, and take $c = 0$. If the starting values are small, such algorithms may take a long time before the output starts to ‘look random’.

There is an interesting observation here. If we define

$$z_n = \sum_{k=0}^{r-1} x_{n-r+k} B^k - \sum_{k=0}^{s-1} x_{n-s+k} B^k + c_{n-1} \quad , \quad (51)$$

then the RCARRY algorithm is seen to be

$$z_n = (a z_{n-1}) \bmod m \quad , \quad m = B^r - B^s + 1 \quad , \quad a = m - (m - 1)/B \quad . \quad (52)$$

It is a multiplicative congruential PRNG, acting on integers that have B digits in base B ; The carry bit is essential to have the multiplications work

out correctly³³. The number $m \approx 2.47 \times 10^{173}$ is prime³⁴, the period is $(m-1)/48$, about 5×10^{171} . In [14] it is recommended to generate numbers in a batch of 24, then skipping the next 223, then another batch of 24, and so on. This is the RANLUX generator.

⌘ A nice feature is that this algorithm can be implemented in floating-point immediately since rounding errors can be avoided.

3.3.6 Good: the Mersenne Twister

Abbreviated by MT, this is likely the most popular PRNG nowadays [15]. It is most convenient to represent the integers x_n in binary form, and then all the operations are in F_2 (that is, addition is binary XOR and multiplication is bitwise AND). We describe here the 32-bit case (hence $w = 32$ below). The algorithm is officially a TGFSRPRNG³⁵, called Mersenne Twister because its period, $2^{19937} - 1$, is a Mersenne prime. The core of the MT algorithm is

$$x_n = x_{n-q} + x_{n-p+1} \begin{pmatrix} 0 & 0 \\ 0 & I_m \end{pmatrix} M + x_{n-p} \begin{pmatrix} I_{w-m} & 0 \\ 0 & 0 \end{pmatrix} M, \quad (53)$$

where I_k is the $k \times k$ unit matrix, and M a carefully chosen $w \times w$ matrix. Note that the binary numbers are multiplied from the left since they are ‘row vectors’. The binary number x_n is then multiplied into a *tempering matrix* T :

$$z_n = x_n T, \quad (54)$$

and z_n is the next pseudorandom integer. The matrix M has a special form:

$$M = \left(\begin{array}{c|c} \begin{matrix} 0 \\ 0 \\ \vdots \end{matrix} & I_{w-1} \\ \hline a \end{array} \right) \quad (55)$$

³³Just like the ‘carry’ operation is necessary in long multiplication.

³⁴For the record:

$m = 247330401473104534060502521019647190035131349101211839914063056092897225106531867170316401061243044987830824361237755009768067533563832694140062258226274209795000570856079361$, and
 $a = 247330386731063812101356613826074608049705993956988322662342632748341364772062482825984947599810524762601263757689206714403985091753014167166773356178267065685142904661606401$.

³⁵Twisted Generalised Feedback Shift Register Pseudo-Random Number Generator.

where a is another w -bit integer. For the tempering step, we can realize that the matrices

$$R = \begin{pmatrix} 0 & & \\ 0 & I_{w-1} & \\ \vdots & & \\ 0 & \dots & 0 \end{pmatrix}, \quad L = \begin{pmatrix} 0 & \dots & 0 \\ & & \vdots \\ I_{w-1} & & 0 \\ & & 0 \end{pmatrix} \quad (56)$$

shift the numbers x to the right and to the left, respectively, by one bit. If we furthermore define the diagonal matrix $D(y) = \text{diag}(y_1, y_2, \dots, y_w)$ for the w -bit number $y = (y_1, y_2, \dots, y_w)$, the tempering matrix T can be written as follows:

$$T = \cdot (R^u \cdot D(d) + 1) \cdot (L^s \cdot D(b) + 1) \cdot (L^t \cdot D(c) + 1) \cdot (R^l + 1) \cdot \quad (57)$$

Here ‘1’ stands for the $w \times w$ unit matrix. In the table we give the values of the various MT parameters.

parameter	decimal	binary
w	32	
p	624	
q	227	
m	31	
a	2567483615	10011001000010001011000011011111
u	11	
d	4294967295	11111111111111111111111111111111
s	7	
b	2636928640	10011101001011000101011010000000
t	15	
c	4022730752	11101111110001100000000000000000
l	18	

The period of MT is $2^{np-w} - 1$ which is huge. Moreover, we define the property of k -distribution to v bits accuracy as follows: taking y_n to consist of the leading v bits of x_n , then the concatenation

$$Y_n = (y_{kn+1}, y_{kn+2}, \dots, y_{kn+k-1})$$

Algorithm 3 The two major steps of the MT algorithm.

{The register is $[x_{n-1}, x_{n-2}, \dots, x_{n-p}]$. Operations: $+$ stands for bitwise XOR, and \times stands for bitwise AND, $>> k$ stands for ‘shift right by k bits’, $<<$ stands for ‘shift left by k bits’.

————— Twisted FSR operation —————

$y \leftarrow (\text{upper } w - r \text{ bits of } x_{n-p+1}) + (\text{lower } r \text{ bits of } x_{n-p})$

if last bit of y is 0 **then**

$y \leftarrow y >> 1$

else

$y \leftarrow (y >> 1) + a$

end if

$x_n \leftarrow x_{n-q} + y$

————— Tempering operation —————

$z \leftarrow x_n + ((x_n >> u) \times d)$

$z \leftarrow z + ((z << s) \times b)$

$z \leftarrow z + ((z << t) \times c)$

$z \leftarrow z + (z >> l)$

return $z/2^{32}$ as a floating-point number

runs over all 2^{kn} possibilities as we progress through the series over its whole period³⁶. Viewed differently, the k -dimensional vectors \vec{Y}_n form a complete and perfect hypercubic lattice³⁷. The MT is k -distributed to 32 bits accuracy for k up to 623, where multiplicative congruential generators typically are no better than 5-distributed.

⌘ However, since in any realistic case the whole period is very far from being used up, the ‘623-distributed’ property is of little use. Indeed the algorithm is known to fail some randomness tests. As one expert states ‘there is no one-size-fits-all PRNG’.

3.4 Exercises

Excercise 4 Randomly chosen algorithms

Perform your own simulation of ‘arbitrary algorithms’. Take M to be largish,

³⁶Except the case where all the y_n are zero: if the register consists of zeroes only, we have a selfie.

³⁷This is either a very good thing or a very bad thing, see also sec.??.

and determine the lifetime for a fixed starting value (1, say). Note that you can do this by simply generating random integers in $(1, M)$ and stopping when you pick some integer for the second time. Do this a great number of times and produce a histogram of the lifetimes, and compute the average.

Excercise 5 RCARRY

Program your own version of the RCARRY algorithm. Experiment with various seeds.

4 Testing PRNGs

4.1 Empirical testing strategies and doubts

Let us consider strings of n bits. Of the 2^n possibilities, some strings, $1111111111\dots$, say, will not ‘look random’, while $10101010101010\dots$ ‘looks slightly more random’. A string like $01101110010111011110001001\dots$ may ‘look quite random’³⁸. In order to decide whether or not to trust a sequence as ‘acceptably random’ it is customary to perform empirical tests on the sequence. That is, one computes *some* number using the sequence, and compares that to what would be expected from an equally long sequence of truly random numbers. Then some criterion is applied to decide whether the string has passed this test. For instance, for $n = 10^4$ we can simply count the number of 1’s. If this falls between, say, 4900 and 5100, we accept the string. Two sobering observations must be made. In the first place, a string of truly random bits will fail this test in about 30 per cent of the cases. We might enlarge our window of acceptance to run from 4800 to 5200 and reject only about 5 per cent of truly random strings, but then we will also accept more nonrandom ones. In the second place, one usually performs several or even many tests. A good string, one that looks appreciable random, will fail about 30 per cent of the tests at the one- σ level: what do you conclude? What if one of the failed tests is actually the integration that you want to perform? When does one stop testing?

4.1.1 The Leeb Conundrum: too much of a good thing

We can look at the procedure of testing n -bit strings in a more abstract manner, adapted from [16]. Let us denote by U the set of all 2^n possible strings. Taking averages for truly random n -bit strings is taking averages over U . A test will single out strings that pass it. We say that the test has *strictness* $1 - s$ if a fraction s of the 2^n strings pass, and all other ones fail. The following sounds trite but is not: *a string x passes test T if it is an element of the subset U_T of U of all strings that pass test T* . In other words, *any* test T , no matter how elaborately formulated, is completely described by U_T . If test T has strictness $1 - s$, then U_T contains $s2^n$ elements. This gives us a handle on the concept of *all tests of strictness $1 - s$* : it is the set of all subsets of U of size $s2^n$. Each subset contains precisely $s2^n$ strings, and

³⁸Actually, it is not! I just wrote 0123456789... in binary.

all strings must occur precisely equally often. The probability that string x passes test T of strictness $1 - s$ is therefore given by

$$\binom{2^n}{s2^n} \times (s2^n) \times \frac{1}{2^n} = \frac{s(2^n)!}{(s2^n)!((1-s)2^n)!} , \quad (58)$$

independently of x . That is, if we perform *all possible tests* of a given strictness, *every* string will perform equally well (or badly) as every other string! Too much testing is not good. In practice, of course, not all tests of a given strictness are considered equivalent, but this only goes to show that in deciding whether or not you like a sequence as ‘random’, art and taste are involved as much as objective testing.

4.1.2 Any test is a uniformity test

4.1.3 The χ^2 characteristic

This quantity is defined for any collection of objects that can be gathered into a finite number of *bins* numbered $1, \dots, B$. Let us suppose that we have some hypothesis (that of equidistribution, for instance) that predicts the *expected* occupancy³⁹ of bin j to be $e_j > 0$. We can compare this to the actual *observed* occupancy of that bin, n_j . The χ^2 statistic tests the observation against the hypothesis as follows:

$$\chi^2 = \sum_{j=1}^B \frac{(n_j - e_j)^2}{e_j} . \quad (59)$$

The Bernoulli distribution assigns the probability of occupancy (n_1, n_2, \dots, n_B) if the objects are independently distributed with probability p_j to end up in bin j :

$$P(n_1, n_2, \dots, n_B) = \frac{N!}{n_1! n_2! \dots n_B!} p_1^{n_1} p_2^{n_2} \dots p_B^{n_B} . \quad (60)$$

Here $N = \sum_j n_j$ is the total number of objects. From

$$\frac{N!}{n_1! n_2! \dots n_j! \dots n_B!} n_j^k = N^k \frac{(N-k)!}{n_1! n_2! \dots (n_j-k)! \dots n_B!} \quad (61)$$

we find the expectations

$$\langle n_j^m \rangle = N^m p_j^m , \quad \langle n_j^{m_j} n_k^{m_k} \rangle = N^{m_j+m_k} p_j^{m_j} p_k^{m_k} . \quad (62)$$

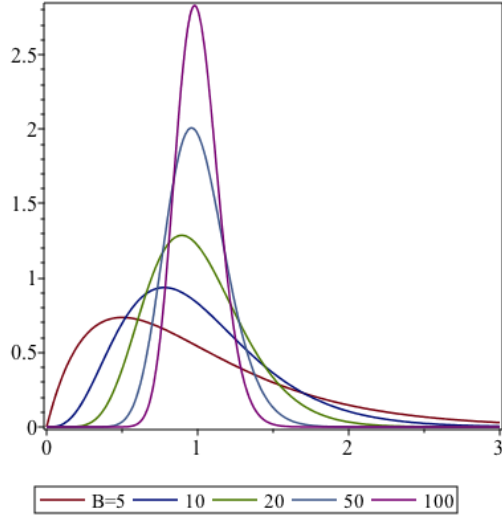
³⁹*i.e.* how many objects fall into the bin.

Some algebra then leads to

$$\begin{aligned}\langle \chi^2 \rangle &= B - 1 , \\ \sigma(\chi^2)^2 &= 2(B - 1) \left(1 - \frac{1}{N} \right) + \frac{1}{N} \left(\sum_{j=1}^B \frac{1}{p_j} - B^2 \right) .\end{aligned}\quad (63)$$

In the variance the second term is nonnegative, vanishing if every $p_j = 1/B$. For large N the density of χ^2 is a Γ -density (*cf* Eq.(147)):

$$P_{\chi^2}(t) = \frac{1}{2\Gamma((B-1)/2)} (t/2)^{(B-3)/2} \exp(-t/2) \quad (64)$$



for χ^2 to attain the value t . The confidence levels (although less simple than those for the normal distribution) are therefore known. The plot shows the renormalised distributions $(B-1)P_{\chi^2}((B-1)t)$ for various values of B . As the number of bins increases the density approaches a Gaussian centered around 1 (*i.e.* the value of χ^2 becomes centered around $B-1$). The nice thing about χ^2 is that for large values of N it only depends on the *number* of bins B . This holds if, say, the second term in Eq.(63) can be neglected. A widely used rule of thumb says that we should make sure that every e_j is at least 5-10 or so. At any rate the *expected* χ^2 is strictly equal to $B-1$, and therefore a quick look at the ‘ χ^2 per degree of freedom’ is usually already a good indicator of a dubious result: it ought to differ not too much from one.

4.2 Theoretical testing strategies

4.2.1 The number-to-number correlation

4.2.2 The spectral test

5 Quasi-Monte Carlo

5.1 Generalities of QMC

5.1.1 The New Leap of Faith

In Monte Carlo integration, the relatively slow $1/\sqrt{N}$ behaviour of the error estimate is due to the fact that the point sets \mathbf{X} are chosen from the ensemble of random iid uniform point sets. We may consider doing better by choosing our point set to be ‘more uniform’ (in some sense to be defined) than a random point set is expected to be. Such point sets are called *superuniform*, or *quasirandom*. Doing Monte Carlo with them is therefore called *Quasi-Monte Carlo* (QMC). The difference with regular Monte Carlo is the fact that the point sets \mathbf{X} used are *not* typical members of the ensemble of random iid uniform point sets⁴⁰. But that means that the *basis* on which the various error estimates $E_{1,2,4}$ of sect. 2.2.2 are constructed is now invalid! We shall have to define a new ensemble, and come to a New Leap of Faith that says that the particular superuniform point set that we use is typical for that ensemble.

It would seem reasonable to insist that, since we take trouble to use point sets with a low measure of nonuniformity, the new superuniform ensemble should be characterised by that measure. For the rest there does not seem to be any reason not to stay as close to the random iid ensemble as possible. Let us assume that there is a function $T(\mathbf{X}) > 0$ that is a measure of the nonuniformity of the point set \mathbf{X} ⁴¹. We also assume that $T(\mathbf{X})$ is invariant under any permutation of the points inside \mathbf{X} . We shall then use the following superuniform ensemble density (restricting ourselves to integration over the unit hypercube I^d):

$$P_S(t; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \frac{1}{P(t)} \delta(T(\mathbf{X}) - t) ,$$

$$P(t) = \int d\mathbf{x}_1 d\mathbf{x}_2 \cdots d\mathbf{x}_N \delta(T(\mathbf{X}) - t) . \quad (65)$$

t is the value of T for all point sets in the superuniform ensemble. The normalisation $P(t)$ is the probability density to find $T(\mathbf{X}) = t$ for a point set \mathbf{X} in the uniform iid ensemble. The density $P(t)$ is computed, for various

⁴⁰Although they are in there, of course; but they form a tiny minority.

⁴¹Hence the smaller $T(\mathbf{X})$, the more uniform the point set \mathbf{X} is.

definitions of nonuniformity T , in section 6. The absence of the iid property is obvious⁴².

⌘ The New Leap of Faith appears to be ignored by almost everyone in the QMC field, with serious consequences (see below).

5.1.2 The mechanism of error improvement

In the superuniform ensemble, the points are not iid as we have seen. On the other hand, by the definition of T and since the superuniform ensemble is a subset of the uniform one, the density $P_S(t; \mathbf{x}_1, \dots, \mathbf{x}_N)$ is symmetric in the point coordinates. Let us consider the situation where we disregard all points except two of them:

$$\frac{1}{P(t)} \int d\mathbf{x}_3 d\mathbf{x}_4 \cdots d\mathbf{x}_N P_S(t; \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N) \equiv 1 - \frac{1}{N} F_2(t; \mathbf{x}_1, \mathbf{x}_2) \quad . \quad (66)$$

This way of introducing the two-point function F_2 is always possible, although we anticipate a factor $1/N$. Moreover we assume that each point in \mathbf{X} is distributed uniformly when considered individually:

$$\int d\mathbf{y} F_2(t; \mathbf{x}, \mathbf{y}) = 0 \quad . \quad (67)$$

If we take the Monte Carlo estimator E_1 to do QMC as well:

$$E_1 = \frac{1}{N} \sum_j w(\mathbf{x}_j) \quad , \quad (68)$$

we see that the superuniform ensemble average

$$\langle E_1 \rangle_S = J_1 \quad (69)$$

again gives us an unbiased integral estimate. The real issue is in the variance:

$$\begin{aligned} \langle E_1^2 \rangle_S &= \frac{1}{N^2} \left\langle \sum_{j,k} w(\mathbf{x}_j) w(\mathbf{x}_k) \right\rangle_S \\ &= \frac{1}{N^2} \left(N \int d\mathbf{x} w(\mathbf{x})^2 + N^2 \left(\int d\mathbf{x} w(\mathbf{x}) \right)^2 \right. \\ &\quad \left. - \frac{N^2}{N} \int d\mathbf{x} d\mathbf{y} F_2(t; \mathbf{x}, \mathbf{y}) w(\mathbf{x}) w(\mathbf{y}) \right) \quad , \quad (70) \end{aligned}$$

⁴²This ensemble may be called *microcanonical*: the corresponding canonical one would let $T(\mathbf{X})$ fluctuate around the value of t , but not much seems to be gained by going over to this ensemble.

so that the variance now reads

$$\sigma(E_1)_S^2 = \frac{1}{N} \left(J_2 - J_1^2 - \left(1 - \frac{1}{N} \right) \int d\mathbf{x} d\mathbf{y} F_2(t; \mathbf{x}, \mathbf{y}) w(\mathbf{x}) w(\mathbf{y}) \right) . \quad (71)$$

Two important conclusions follow:

1. A small, $\mathcal{O}(1/N)$ deviation from iid uniformity can have a large effect on the expected integration error.
2. We may expect that the integrand values $w(\mathbf{x})$ and $w(\mathbf{y})$ are correlated when \mathbf{x} and \mathbf{y} are close to one another⁴³; in particular they will tend to have the same sign. For QMC to be an improvement, we therefore want $F_2(t; \mathbf{x}, \mathbf{y})$ to become positive when \mathbf{x} and \mathbf{y} approach one another: in a sense, the points in \mathbf{X} must feel a mutual repulsion.

Another way to write Eq.(71) is

$$\sigma(E_1)_S^2 = \frac{1}{2N} \int d\mathbf{x} d\mathbf{y} \left(1 + \frac{N-1}{N} F_2(t; \mathbf{x}, \mathbf{y}) \right) \left(w(\mathbf{x}) - w(\mathbf{y}) \right)^2 . \quad (72)$$

The ideally best possible error estimate that is valid for any integrand w is therefore reached for

$$F_2(t; \mathbf{x}, \mathbf{y}) = -1 + \delta^d(\mathbf{x} - \mathbf{y}) . \quad (73)$$

As we shall see, this is precisely what we obtain for $T(\mathbf{X}) \rightarrow 0$, a beautiful but admittedly unreachable situation.

We of course have to somehow compute $P(t)$ and F_2 for a given nonuniformity T , and show that F_2 integrates to zero. This is dealt with in section 6.

⌘ Like the New Leap of Faith, the mechanism behind the error improvement of QMC over MC seems to be virtually unknown.

5.2 Error estimators

5.2.1 The first-order estimate

5.2.2 The second-order estimate

5.2.3 Payback time: Lack of Leap of Faith is Punished

⁴³In some reasonable sense.

6 Nonuniformity of point sets

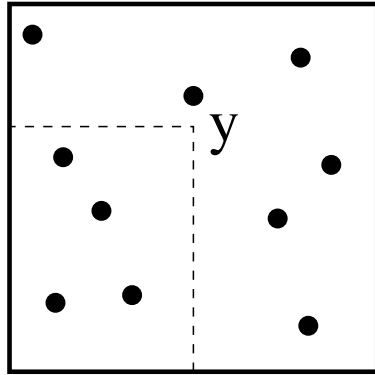
6.1 Measures of nonuniformity: Discrepancy

6.1.1 The star discrepancy

Consider *any* point set $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ in the d -dimensional unit hypercube I^d . With $\theta(\mathbf{x} < \mathbf{y})$ we shall mean the function that is 1 if $\mathbf{x}^\mu < \mathbf{y}^\mu$ for all $\mu = 1, 2, \dots, d$, otherwise zero⁴⁴. We define the *local discrepancy* $g(\mathbf{y})$ as

$$g(\mathbf{y}) = \frac{1}{N} \sum_j h(\mathbf{x}_j; \mathbf{y}) ,$$

$$h(\mathbf{x}_j; \mathbf{y}) = \theta(\mathbf{x}_j < \mathbf{y}) - \text{vol}(\mathbf{y}) , \quad \text{vol}(\mathbf{y}) = \prod_{\mu=1}^d \mathbf{y}^\mu . \quad (74)$$



The local discrepancy compares the fraction of points ‘below’ \mathbf{y} with what that fraction would be if the points were ideally uniformly distributed. In the plot we give an example for $d = 2$. The point \mathbf{y} has coordinates $(1/2, 2/3)$ and its rectangle contains 4 points out of 10, hence $g(\mathbf{y}) = 1/15$. At every corner, the local discrepancy vanishes. Measures of *global*

discrepancy can be defined: the most important are the *extreme* discrepancy⁴⁵ (the *Kolmogorov-Smirnov* statistic):

$$L_\infty^* = \sup_{\mathbf{y} \in I^d} |g(\mathbf{y})| , \quad (75)$$

and the *quadratic* discrepancy (the *Kramér-von Mises* statistic):

$$L_2^* = \int_{I^d} d\mathbf{y} g(\mathbf{y})^2 . \quad (76)$$

⁴⁴That is, \mathbf{x} finds itself inside the rectangle spanned by \mathbf{y} and the origin.

⁴⁵The asterisk refers to the fact that the hyper-rectangles are attached to the origin.

The extreme discrepancy is beloved of mathematicians, but the quadratic one is easier to handle in computations. At any rate, if one is small the other will be small as well, since $L_2^* < (L_\infty^*)^2$ and the function $g(\mathbf{y})$ is piecewise linear in the components \mathbf{y}^μ ⁴⁶. The L_2^* discrepancy can be formulated as a function of \mathbf{X} only:

$$L_2^* = \frac{1}{N^2} \sum_{i,j=1}^N \prod_{\mu=1}^d \left(1 - \max(\mathbf{x}_i^\mu, \mathbf{x}_j^\mu)\right) - \frac{2}{2^d N} \sum_{i=1}^N \prod_{\mu=1}^d \left(1 - (\mathbf{x}_i^\mu)^2\right) + \left(\frac{1}{3}\right)^d . \quad (77)$$

For point sets \mathbf{X} taken from the random iid ensemble, we have the expectation

$$\langle L_2^* \rangle = \frac{1}{N} \left(2^{-d} - 3^{-d}\right) . \quad (78)$$

It is instructive to compare this for the ‘hypercubic’ lattice with $N = M^d$ points. This has

$$x_j^\mu = \frac{2k_\mu - 1}{2M} , \quad j = 1 + \sum_{\mu=1}^D (k_\mu - 1) M^{\mu-1} , \quad k_\mu \in \{1, 2, \dots, M\} , \quad (79)$$

and the quadratic discrepancy evaluates to

$$L_2^* = \left(\frac{1}{3}\right)^d \left[\left(1 + \frac{1}{2M^2}\right)^d - 2 \left(1 + \frac{1}{8M^2}\right)^d + 1 \right] \approx \frac{d}{4N^{2/d} 3^d} , \quad (80)$$

where the approximation holds for large M . For $d > 4$ the ‘random’ point sets are, by this definition, *more uniform* than the regular hypercubic ones.

⋈ For given N , the hypercubic lattice will win out as d increases. However, the usual situation is that d is fixed (by the integration problem itself), and N is the free parameter that tells us how much computing resources can be spent.

6.1.2 Random *vs* Regular: Translation *vs* Rotation

You may wonder how a random point set can be more uniform than a nice, beautiful hypercubic lattice. The answer⁴⁷ must be the following. In low

⁴⁶With jumps of magnitude $1/N$ whenever $\mathbf{y}^\mu = \mathbf{x}_j^\mu$ for some μ and j .

⁴⁷At least, the answer that satisfies me.

dimensions ($d = 1, 2$ or so) the translational invariance (by steps of $1/M$) of the hypercubic lattice is an obvious advantage. Indeed, in one dimension the regularly-spaced point set is the most uniform one by any standard. However, in more dimensions we also have to consider the *rotational* properties of the point set. A random collection of points does not look very different when rotated, whereas for the hypercubic lattice some directions contain many points (especially if we look parallel to the axes), while some directions contain hardly any points at all. As the dimensionality increases, there are more and more directions to choose from, and the rotational invariance becomes the more dominant property.

6.1.3 The Roth bound

It is obvious that the ideal finite point set, with $L_2^* = L_\infty^* = 0$, does not exist. In fact the discrepancies have lower bounds. In [17]⁴⁸ the so-called *Roth bound* is proven: for any point set $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ in I^d ($d \geq 2$) there is a constant $c_d > 0$ (depending only on d) such that

$$L_2^*(\mathbf{X}) > c_d \frac{\log(N)^{d-1}}{N^2} . \quad (81)$$

We see that the discrepancy can be quite a bit smaller than the ‘expected’ $(2^{-d} - 3^{-d})/N$ for random point sets; the question, of course, is how to find such low-discrepancy point sets, and what c_d is. In [17] it is proven that Eq.(81) holds with

$$c_d = 2^{-8d} \left((d-1) \log(2) \right)^{1-d} , \quad (82)$$

but I believe the ‘real’ c_d is (considerably) larger.

6.1.4 The Koksma-Hlawka inequality

6.1.5 The Wiener measure and the Woźniakowski Lemma

Let us consider functions $f(\mathbf{x})$ on the d -dimensional unit hypercube I^d . For $d = 1$, the *Wiener measure* W describes an ensemble of functions that is

⁴⁸This book is an absolute must for anyone interested in low-discrepancy point sets and related matters. I find it interesting to see that the emphasis is on L_∞^* rather than on L_2^* , perhaps reflecting the background of the authors as mathematicians rather than physicists. Section 2.2, Lemma 2.5 is especially relevant.

defined by its expectation values:

$$\langle f(x) \rangle_W = 0 \quad , \quad \langle f'(x) f'(y) \rangle_W = \delta(x - y) \quad . \quad (83)$$

By integration we then find

$$\langle f(x) f(y) \rangle_W = \min(x, y) \quad . \quad (84)$$

Let us now consider integrating such a function using Monte Carlo. The integration error,

$$\eta = \frac{1}{N} \sum_j f(x_j) - \int_0^1 dx f(x) \quad , \quad (85)$$

(where the sum runs from $j = 1$ to $j = N$) has a squared average over W , called the *complexity*:

$$\langle \eta^2 \rangle_W = \frac{1}{N^2} \sum_{j,k} \min(x_j, x_k) - \frac{2}{N} \sum_j (x_j - x_j^2/2) + \frac{1}{3} \quad . \quad (86)$$

The more-dimensional variant is called the *Wiener sheet measure*:

$$\langle \varphi(\mathbf{x}) \rangle_W = 0 \quad , \quad \left\langle \frac{\partial^n}{\partial x^1 \dots \partial x^d} f(\mathbf{x}) \frac{\partial^n}{\partial y^1 \dots \partial y^d} f(\mathbf{y}) \right\rangle_W = \delta^d(\mathbf{x} - \mathbf{y}) \quad , \quad (87)$$

and

$$\langle f(\mathbf{x}) f(\mathbf{y}) \rangle_W = \prod_{\mu=1}^d \min(\mathbf{x}^\mu, \mathbf{y}^\mu) \quad . \quad (88)$$

The corresponding complexity is

$$\langle \eta^2 \rangle_W = \frac{1}{N^2} \sum_{j,k} \prod_{\mu} \min(\mathbf{x}_j^\mu, \mathbf{x}_k^\mu) - \frac{2}{N} \sum_j \left(\mathbf{x}_j^\mu - \frac{1}{2} (\mathbf{x}_j^\mu)^2 \right) + \left(\frac{1}{3} \right)^d \quad . \quad (89)$$

By replacing every \mathbf{x}^μ by $1 - \mathbf{x}^\mu$ this is seen to be nothing but the quadratic discrepancy L_2^* , with every (hyper)rectangle anchored not to the origin point $(0, 0, \dots, 0)$ but the point $(1, 1, \dots, 1)$. This is the Woźniakowski lemma: for functions drawn from the Wiener (sheet) measure the complexity is given by the discrepancy of the point set [18]. The central notion here is that of a *problem class*: the ensemble of functions that we assume our particular integrand is drawn from.

⌘ The type of function typical of the Wiener (sheet) measure is hardly what you would encounter in high-energy phenomenology: functions with a fractal structure, continuous but nowhere differentiable. Our usual integrand, in contrast, has discontinuities but is differentiable elsewhere. Nevertheless the above derivation shows that a conjectured ensemble of integrands dictates a measure of nonuniformity. In what follows we try to make this notion more workable.

6.2 Measures of nonuniformity: Diaphony

6.2.1 Fourier problem classes

We consider the set of all functions f periodic on the d -dimensional unit hypercube that can be written using Fourier modes,

$$f(\mathbf{x}) = \sum_{\mathbf{n}}^* a_{\mathbf{n}} \exp(2i\pi \mathbf{n} \cdot \mathbf{x}) , \quad (90)$$

where the asterisk means that the sum is to be taken over all d -dimensional vectors \mathbf{n} with integer components, except the zero vector $\mathbf{n} = (0, 0, \dots, 0)$. Since constant functions are integrated with zero error, the zero modes are irrelevant to our discussion⁴⁹. We take the $a_{\mathbf{n}}$ to be real numbers. The Fourier ensemble measure is then defined by the measure on the set of coefficients $a_{\mathbf{n}}$, which we take to be Gaussian as follows:

$$\mathcal{D}f = \prod_{\mathbf{n}}^* da_{\mathbf{n}} \exp\left(-\frac{a_{\mathbf{n}}^2}{2\sigma_{\mathbf{n}}^2}\right) \frac{1}{\sqrt{2\pi\sigma_{\mathbf{n}}^2}} \quad (91)$$

This implies the following ensemble averages:

$$\langle a_{\mathbf{n}} \rangle_F = 0 \quad , \quad \langle a_{\mathbf{n}} a_{\mathbf{n}'} \rangle_F = \sigma_{\mathbf{n}}^2 \theta(\mathbf{n} = \mathbf{n}') . \quad (92)$$

The real quantity $\sigma_{\mathbf{n}}^2$ is called the *strength* of the mode \mathbf{n} . We shall assume that $\sigma_{\mathbf{n}}^2$ depends on the components \mathbf{n}^μ ($\mu = 1, 2, \dots, d$) only through their absolute values. Moreover we shall insist that the total strength is finite, in fact by convention

$$\sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 = 1 . \quad (93)$$

⁴⁹Indeed, you may argue that any integration is nothing but an attempt to strip an integrand of all its *nonzero* modes.

⌘ The functions f are real *on average*. The Fourier problem class contains the real-valued functions but is actually even larger.

6.2.2 Fourier diaphony

We envisage integrating f using a point set $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, so that the integration error is the estimated integral:

$$\eta = \frac{1}{N} \sum_{j=1}^N f(\mathbf{x}_j) \quad (94)$$

The squared error, averaged over the Fourier problem class, is then

$$\begin{aligned} \langle |\eta|^2 \rangle_F &= \frac{1}{N^2} \sum_{j,k=1}^N \sum_{\mathbf{n}, \mathbf{n}'}^* \langle a_{\mathbf{n}} a_{\mathbf{n}'} \rangle \exp \left(2i\pi (\mathbf{n} \cdot \mathbf{x}_j - \mathbf{n}' \cdot \mathbf{x}_k) \right) \\ &= \frac{1}{N^2} \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 \left| \sum_{j=1}^N \exp(2i\pi \mathbf{n} \cdot \mathbf{x}_j) \right|^2 . \end{aligned} \quad (95)$$

The better the point set is at integrating the various Fourier modes, the smaller will be the error. This, then, leads us to define the *diaphony* of the point set as

$$\begin{aligned} T(\mathbf{X}) &= \frac{1}{N} \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 \left| \sum_{j=1}^N \exp(2i\pi \mathbf{n} \cdot \mathbf{x}_j) \right|^2 = \frac{1}{N} \sum_{j,k=1}^N \beta(\mathbf{x}_j - \mathbf{x}_k) , \\ \beta(\mathbf{z}) &= \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 \exp(2i\pi \mathbf{n} \cdot \mathbf{z}) . \end{aligned} \quad (96)$$

The *bare two-point function* $\beta(\mathbf{z})$ is periodic in each of the components \mathbf{z}^μ , and is invariant under $\mathbf{z}^\mu \rightarrow -\mathbf{z}^\mu$ for $\mu = 1, 2, \dots, \mu$. In addition we have the following important properties:

$$\beta(0) = 1 \quad , \quad \int \beta(\mathbf{z}) d\mathbf{z} = 0 . \quad (97)$$

The factor in front of the definition of $T(\mathbf{X})$ reads $1/N$ (not $1/N^2$) by convention. The choice of the strengths $\sigma_{\mathbf{n}}^2$ specifies the particular diaphony. By construction, $T(\mathbf{X})$ is nonnegative, and for the ‘most nonuniform’ point set, where all \mathbf{x}_j coincide, $T(\mathbf{X}) = N$.

⌘ To have a diaphony that we can realistically compute for a point set, it is important to have $\beta(\mathbf{z})$ in some kind of closed form.

6.2.3 Choosing your strengths: examples of diaphony

A few examples of diaphony may be given. These are mainly geared towards allowing results in closed form, plus the ‘physical’ intuition that modes with small $|\mathbf{n}|$ are ‘naturally’ more prominent than those with large $|\mathbf{n}|$.

1. Euler diaphony⁵⁰ T_E for $d = 1$:

$$\begin{aligned}\sigma_n^2 &= \frac{3}{\pi^2} \frac{1}{n^2} , \\ \beta_E(x) &= 1 - 6|x|(1 - |x|)\end{aligned}\tag{98}$$

For the computation of $\beta_E(x)$ see sec.11.0.7.

2. Euler diaphony T_E for $d > 1$:

$$\begin{aligned}\sigma_{\mathbf{n}}^2 &= \frac{1}{2^d - 1} \prod_{\mu=1}^d \left(\theta(\mathbf{n}^\mu = 0) + \frac{3}{\pi^2 (\mathbf{n}^\mu)^2} \theta(\mathbf{n}^\mu \neq 0) \right) , \\ \beta(\mathbf{x}) &= \frac{1}{2^d - 1} \left(-1 + \prod_{\mu=1}^d (1 + \beta_E(\mathbf{x}^\mu)) \right)\end{aligned}\tag{99}$$

3. Gulliver diaphony⁵¹ T_G :

$$\begin{aligned}\sigma_{\mathbf{n}}^2 &= \left(\left(\frac{1+s}{1-s} \right)^d - 1 \right)^{-1} s^{-|n^1| - |n^2| - \dots - |n^d|} , \quad 0 < s < 1 , \\ \beta_G(\mathbf{x}) &= \left(\left(\frac{1+s}{1-s} \right)^d - 1 \right)^{-1} \left(-1 + \prod_{\mu=1}^d \frac{1 - s^2}{1 - 2s \cos(2\pi \mathbf{x}^\mu) + s^2} \right)\end{aligned}\tag{100}$$

4. Block diaphony T_B :

$$\begin{aligned}\sigma_{\mathbf{n}}^2 &= \frac{1}{2p} \prod_{\mu=1}^d \theta(-c \leq \mathbf{n}^\mu \leq c) , \quad 2p = (2c + 1)^d - 1, \\ \beta_B(\mathbf{x}) &= \frac{1}{2p} \left(-1 + \prod_{\mu=1}^d \frac{\sin((2c + 1)\pi \mathbf{x}^\mu)}{\sin(\pi \mathbf{x}^\mu)} \right)\end{aligned}\tag{101}$$

⁵⁰Named after Euler’s formula $\sum_{n \geq 1} n^{-2} = \pi^2/6$.

⁵¹After Gulliver de Boer, the student who suggested it.

5. Jacobi diaphony⁵² T_J :

$$\begin{aligned}
\sigma_{\mathbf{n}}^2 &= \frac{1}{K(\lambda; 0)^d - 1} \exp(-\lambda |\mathbf{n}|^2) , \\
\beta_J(\mathbf{x}) &= \frac{1}{K(\lambda; 0)^d - 1} \left(-1 + \prod_{\mu=1}^d K(\lambda; \mathbf{x}^\mu) \right) \\
K(\lambda; z) &= \sum_n \exp(-\lambda n^2 + 2i\pi n z) \\
&= \sqrt{\frac{\pi}{\lambda}} \sum_n \exp\left(-\frac{\pi^2(n+z)^2}{\lambda}\right) .
\end{aligned} \tag{102}$$

In the last line the Poisson summation formula has been used (*cf* sec.11.0.3).

A major drawback of diaphonies is that lattice vectors \mathbf{n} that have the same norm but are oriented differently can have very different strengths, especially the more-dimensional Euler diaphony. The Gulliver diaphony aims at repairing this somewhat while still having a closed form for the two-point function. The Jacobi diaphony comes closest to full rotational invariance, while the sums can usually be restricted to a manageable number of terms.

⌘ For the Euler diaphony in two dimensions, the lattice vectors $\mathbf{n} = (13, 0)$ and $\mathbf{n} = (12, 5)$ make an angle of only 21 degrees, but their strengths differ by a factor 21 as well.

6.3 QFT for diaphony

6.3.1 The distribution of diaphony

A given point set may have a diaphony of, say, 0.8: is this good, or bad, or what? The natural yardstick is of course what you would expect for a *random* point set. It is therefore interesting to see what can be said about the distribution of $T(\mathbf{X})$ if the point set \mathbf{X} is taken from the ensemble of random point sets discussed in section 2.1.3. In particular, we will be interested in $\langle T^\ell \rangle$ ($\ell = 1, 2, \dots$) where the average is over the ensemble of random point sets, that is, we integrate over the points in the point set, assuming them to be iid uniform in the hypercube. Eventually, we want to be able to say something about the generating function $\Omega(z) \equiv \langle \exp(zT) \rangle$.

⁵²Beause of the emergence of Jacobi theta functions with real-valued nome.

6.3.2 Feynman rules for diaphony in the large- N limit

It is useful to formulate the computation of the various moments of T in terms of Feynman diagrams. The bare two-point function $\beta(\mathbf{x}_j - \mathbf{x}_k)$ is the bare propagator, and the points themselves are the vertices of the diagrams. The lattice vectors \mathbf{n} are the momenta, and momentum is conserved at each vertex. An example is

$$\text{---} \bullet \text{---} \bigcirc \text{---} \bullet \text{---} = \int d\mathbf{x}_1 d\mathbf{x}_2 d\mathbf{x}_3 d\mathbf{x}_4 \beta(\mathbf{x}_1 - \mathbf{x}_2) \beta(\mathbf{x}_2 - \mathbf{x}_3)^2 \beta(\mathbf{x}_3 - \mathbf{x}_4)$$

Diagrams can be disconnected. If a diagram contains k vertices then it picks up a factor N^k since we then have to sum over k *distinct* points. In addition there is a combinatorial factor, the number of ways that diagram can be formed. The first moments of the T distribution are now

$$\begin{aligned} \langle T \rangle &= \frac{1}{N} \left(N^2 \text{---} \bullet \text{---} \bullet \text{---} + N^1 \bigcirc \right) , \\ \langle T^2 \rangle &= \frac{1}{N^2} \left(N^4 \begin{array}{c} \bullet \text{---} \bullet \\ \bullet \text{---} \bullet \end{array} + 2N^3 \begin{array}{c} \bullet \text{---} \bullet \\ \bullet \text{---} \bullet \end{array} \bigcirc + 4N^3 \text{---} \bullet \text{---} \bullet \text{---} \bullet \text{---} + 4N^2 \text{---} \bullet \text{---} \bigcirc \\ &\quad + 2N^2 \bigcirc \bullet \text{---} \bullet \text{---} \bigcirc + N^2 \bigcirc \bullet \text{---} \bullet \text{---} \bigcirc + N^1 \bigcirc \bullet \text{---} \bullet \text{---} \bigcirc \right) \end{aligned} \quad (103)$$

A tremendous simplification arises from the properties of $\beta(\mathbf{z})$. In the first place, all tadpoles vanish:

$$\text{---} \bullet \text{---} \bigcirc = 0 . \quad (104)$$

Secondly, we have *one-vertex reducibility*: any two pieces of a connected diagram that are connected by a single vertex factor can be separated⁵³:

$$\text{---} \bullet \text{---} \bigcirc \text{---} \bullet \text{---} = \text{---} \bullet \text{---} \bigcirc \text{---} \bullet \text{---} \text{---} \bullet \text{---} \bigcirc \text{---} \bullet \text{---} . \quad (105)$$

This is due to the fact that no momentum can flow between the pieces through the connecting vertex. All this means that only one-particle irreducible vacuum diagrams are nonzero. In addition, we are not really interested in values

⁵³Without changing the N^k in front, of course.

of N smaller than 10^3 or so. We can therefore take the large- N limit, so that in $\langle T^p \rangle$ only those diagrams survive that carry a prefactor $N^p \approx N^p$:

$$\begin{aligned}\langle T \rangle &= \text{diagram with 1 bead} , \\ \langle T^2 \rangle &= 2 \text{diagram with 2 beads} + \text{diagram with 2 beads}^2 , \\ \langle T^3 \rangle &= 8 \text{diagram with 3 beads} + 6 \text{diagram with 2 beads} \text{diagram with 2 beads} + \text{diagram with 3 beads}^3 ,\end{aligned}\tag{106}$$

and so on⁵⁴. These diagrams are sums of products of *bracelets*. A k -bead bracelet B_k evaluates to

$$B_k \equiv \int d\mathbf{x}_1 d\mathbf{x}_2 \cdots d\mathbf{x}_k \beta(\mathbf{x}_1 - \mathbf{x}_2) \beta(\mathbf{x}_2 - \mathbf{x}_3) \cdots \beta(\mathbf{x}_k - \mathbf{x}_1) = \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^{2k} . \tag{107}$$

Owing to our choice of diaphony⁵⁵ the diaphony distribution becomes N -independent for large N . In other terms in the $1/N$ expansion we may encounter non-bracelet diagrams such as the last diagram in Eq.(103). In general, a L -loop connected diagram carries a factor $1/N^{L-1}$, and we see that $1/N$ plays the rôle of Planck's constant.

⌘ The star discrepancy $N L_2^*$ of course has its own propagator, defined as

$$\beta(\mathbf{x}_1, \mathbf{x}_2) = \int d\mathbf{y} h(\mathbf{x}_1; \mathbf{y}) h(\mathbf{x}_2; \mathbf{y})$$

But since this propagator is not translationally invariant nor integrates to zero, its analysis in terms of Feynman diagrams is much more complicated. Still, some results have been derived [19].

6.3.3 Collecting bracelets

We can immediately find the expectation value⁵⁶ and variance of T :

$$\langle T \rangle = \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 = 1 \quad , \quad \sigma(T)^2 = \langle T^2 \rangle - 1 = 2 \sum_{\mathbf{n}} \sigma_{\mathbf{n}}^4 . \tag{108}$$

⁵⁴The sum of the coefficients in $\langle T^m \rangle$ is $(2m)!/2^m m!$.

⁵⁵In particular the use of the factor $1/N$ rather than $1/N^2$.

⁵⁶This was the reason for insisting on the normalisation of the total strength.

But we can do more. Let us consider a contribution to $\langle \exp(zT) \rangle$ that contains n_1 one-bead bracelets, n_2 two-bead bracelets, n_3 three-bead bracelets, and so on. The total number of propagators involved is then

$$n = n_1 + 2n_2 + 3n_3 + \dots \quad (109)$$

The number of ways to divide n propagators into n_1 groups of one, n_2 groups of two, n_3 groups of three, and so on, is

$$R(n_1, n_2, n_3, \dots) = \frac{n!}{n_1! n_2! n_3! \dots (1!)^{n_1} (2!)^{n_2} (3!)^{n_3} \dots} \quad , \quad (110)$$

where we also take into account the indistinguishability of the propagators in each group. Now, a k -bead bracelet can be built from k propagators in

$$2(k-1)2(k-2)2(k-3)\dots 2 = \frac{2^k k!}{2k}$$

ways. The contribution under consideration has, in addition, a factor $z^n/n!$. Putting everything together and summing over all values of n_1, n_2, n_3, \dots then gives us

$$\begin{aligned} \Omega(z) &= \sum_{n_1, n_2, n_3, \dots \geq 0} \frac{1}{n_1!} \left(\frac{2z}{2} B_1 \right)^{n_1} \frac{1}{n_2!} \left(\frac{(2z)^2}{4} B_2 \right)^{n_2} \frac{1}{n_3!} \left(\frac{(2z)^3}{2} B_3 \right)^{n_3} \dots \\ &= \exp \left(\sum_{k \geq 1} \frac{(2z)^k}{2k} B_k \right) = \exp \left(-\frac{1}{2} \sum_{\mathbf{n}}^* \log(1 - 2z\sigma_{\mathbf{n}}^2) \right) \\ &= \left(\prod_{\mathbf{n}}^* (1 - 2z\sigma_{\mathbf{n}}^2) \right)^{-1/2} . \end{aligned} \quad (111)$$

In this derivation, the factor $1/(2k)$ is precisely the symmetry factor of the k -bead bracelet, and the occurrence of the 2 in the $(2z)^k$ is due to the bosonic nature of the propagators, by which they can always be connected in two ways in any bracelet. For the block diaphony we have

$$\Omega_B(z) = \frac{1}{(1 - z/p)^p} \quad , \quad (112)$$

and for the one-dimensional Euler diaphony

$$\Omega_E(z) = \prod_{n \geq 1} \left(1 - \frac{6z}{\pi^2 n^2} \right)^{-1} = \frac{\sqrt{6z}}{\sin(\sqrt{6z})} \quad . \quad (113)$$

The one-dimensional Gulliver diaphony has

$$\Omega_G(z) = \prod_{n \geq 1} (1 - s^n x)^{-1} \ , \ x = \frac{1-s}{s} z \ . \quad (114)$$

6.3.4 The diaphony distribution for large N

The actual probability density $P(t) = \text{Prob}(T(\mathbf{X}) = t)$ is given by the inverse Laplace transform

$$\begin{aligned} P(t) &= \frac{1}{2\pi i} \int_{\Gamma} dz \exp(-zt) \Omega(z) \\ &= \frac{1}{2\pi} \int_{\Gamma} dz \exp\left(-zt - \frac{1}{2} \sum_{\mathbf{n}}^* \log(1 - 2z\sigma_{\mathbf{n}}^2)\right) \ , \end{aligned} \quad (115)$$

where the integration contour Γ runs from $-i\infty$ to $+i\infty$, passing to the left of every singularity of $\Omega(z)$, that is it crosses the real axis below $1/(2 \max_{\mathbf{n}} \sigma_{\mathbf{n}}^2)$. For some diaphonies we can write it in (almost) closed form: for the block diaphony

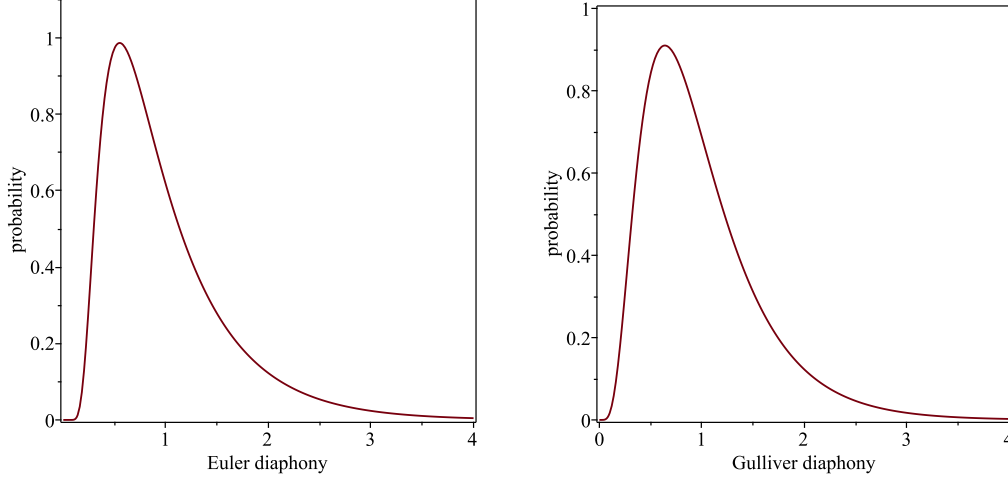
$$P_B(t) = \frac{p^p}{(p-1)!} t^{p-1} \exp(-tp) \ , \ p = ((2c+1)^d - 1)/2 \ , \quad (116)$$

and for the Euler diaphony (again employing Poisson's formula)

$$\begin{aligned} P_E(t) &= \sum_{n \geq 1} (-)^{n-1} \frac{n^2 \pi^2}{3} \exp\left(-\frac{n^2 \pi^2 t}{6}\right) \\ &= \sqrt{\frac{3}{2\pi t}} \sum_{n=-\infty}^{\infty} \left(\frac{3(2n+1)^2}{t^2} - \frac{1}{t}\right) \exp\left(-\frac{3(2n+1)^2}{2t}\right) \ . \end{aligned} \quad (117)$$

Finally, for the one-dimensional Gulliver diaphony we find

$$\begin{aligned} P_G(t) &= \sum_{m \geq 1} \exp\left(\frac{-t}{(1-s)s^{m-1}}\right) R_m(s) \ , \\ R_1(s) &= \frac{1}{(1-s)} \prod_{n \geq 1} (1 - s^n)^{-1} \ , \\ R_m(s) &= -\frac{s^{m-2}}{1 - s^{m-1}} R_{m-1}(s) \ , \ m \geq 2 \ . \end{aligned} \quad (118)$$



The plots show the large- N probability density of the one-dimensional Euler and Gulliver ($s = 0.5$) diaphonies. For large t the distribution decays exponentially. The maximum probability is attained for t values considerably lower than the mean value 1. At $t = 0$, the Euler diaphony distribution $P_E(t)$ has an essential singularity. For all diaphonies that have an infinite number of nonzero strengths the distribution $P(t)$ must have the form

$$P(t) = \sum_n c_n \exp(-a_n t) \quad , \quad \sum_n c_n = 0 \quad ,$$

where the c_n asymptotically go to zero while the a_n increase without bound⁵⁷. As soon as $\Re(t) < 0$ the exponentials will explode and $P(t)$ is no longer finite⁵⁸. We therefore conjecture that $t = 0$ is a singular point of $P(t)$ if the number of nonzero strengths is infinite.

6.3.5 The saddle-point approximation

In the computation of $P(t)$ we may use a saddle-point approximation, as follows. We can write

$$P(t) = \frac{1}{2\pi i} \int_{\Gamma} dz \exp(\phi_t(z)) \quad , \quad \phi_t(z) = -zt - \frac{1}{2} \sum_{\mathbf{n}}^* \log(1 - 2z\sigma_{\mathbf{n}}^2) \quad . \quad (119)$$

⁵⁷Otherwise the total strength would not be finite.

⁵⁸I have checked this numerically for the Gulliver diaphony.

We can look for the extremal point z_0 of $\phi_t(z)$:

$$\phi'_t(z_0) = \sum_{\mathbf{n}}^* \frac{\sigma_{\mathbf{n}}^2}{1 - 2z_0\sigma_{\mathbf{n}}^2} - t = 0 \quad , \quad z_0 < 1/(2 \max_{\mathbf{n}} \sigma_{\mathbf{n}}^2) \quad . \quad (120)$$

Then we can write the saddle-point approximation as

$$\begin{aligned} \phi_t(z) &\approx \phi_t(z_0) + \frac{1}{2}(z - z_0)^2 \phi''_t(z_0) \quad , \quad \phi''_t(z_0) = \sum_{\mathbf{n}}^* \frac{2\sigma_{\mathbf{n}}^4}{(1 - 2z_0\sigma_{\mathbf{n}}^2)^2} \quad , \\ P(t) &\approx \sqrt{\frac{2\pi}{\phi''_t(z_0)}} \exp(\phi_t(z_0)) \quad . \end{aligned} \quad (121)$$

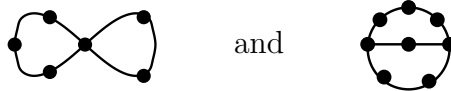
For the block diaphony this is simply the Stirling approximation for the prefactor $p^p/(p-1)!$, good to better than one percent even at $p \approx 10$. For the Euler and Gulliver diaphonies the approximation is also excellent. In practice, to obtain the saddle-point results it is best to take z_0 as the independent variable, and compute both t and $P(t)$ as a function of z_0 between $-\infty$ and the first singularity. Note that $t \rightarrow \infty$ corresponds to having z_0 sidling up to the first singularity. The saddle-point value $z_0 = 0$ corresponds, very properly, to the expectation value $t = 1$. The limit $t \downarrow 0$ is reflected in $z_0 \rightarrow -\infty$. For the Block diaphony we have the exact relation $z_0 = p(1 - 1/t)$.

6.3.6 $1/N$ corrections to the diaphony distribution

So far we have discussed only the $N \rightarrow \infty$ limit for $P(t)$. In order to obtain the leading correction in the $1/N$ expansion we must take into account the following two effects. In the first place, concomitant with every factor z^q the product of all bracelets contains the factor $N^q/N^q \approx 1 - q(q-1)/2N$. This $1/N$ effect can be represented by taking

$$\left(1 - \frac{z^2}{2N} \frac{\partial^2}{(\partial z)^2}\right) \Omega(z) = \Omega(z) \left[1 - \frac{z^2}{2N} \left(\left(\sum_{\mathbf{n}}^* \frac{\sigma_{\mathbf{n}}^2}{1 - 2z\sigma_{\mathbf{n}}^2} \right)^2 + \sum_{\mathbf{n}}^* \frac{2\sigma_{\mathbf{n}}^4}{1 - 2z\sigma_{\mathbf{n}}^2} \right) \right] \quad . \quad (122)$$

In the second place, we must include diagrams that have one vertex less than the number of its propagators⁵⁹: these are of the forms



⁵⁹But of these diagrams we need to take only the first power.

with any number of vertices added onto the various propagators. It is seen that p -point vertices effectively carry a coupling constant $N^{-(p-2)/2}$. It is useful to introduced the *dressed propagator*

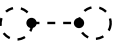
$$\times \text{-----} \times \equiv \tilde{\beta}(z; \mathbf{x}) = \sum_{\mathbf{n}}^* \frac{2z\sigma_{\mathbf{n}}^2}{1 - 2z\sigma_{\mathbf{n}}^2} \exp(2i\pi \mathbf{n} \cdot \mathbf{x}) \quad , \quad (123)$$

denoted by a dashed line, in terms of which we can write the generating function including its $1/N$ corrections as

$$\langle e^{zT} \rangle = \Omega(z) \left(1 - \frac{1}{8N} \text{[diagram 1]}^2 - \frac{1}{4N} \text{[diagram 2]} + \frac{1}{8N} \text{[diagram 3]} + \frac{1}{12N} \text{[diagram 4]} \right) \quad (124)$$

Note that all these diagrams carry their ‘natural’ symmetry factor. The first and third of the four diagrams cancel one another and we find

$$\langle e^{zT} \rangle = \Omega(z) \left(1 - \frac{1}{4N} \int d\mathbf{x} \tilde{\beta}(z; \mathbf{x})^2 + \frac{1}{12N} \int d\mathbf{x} \tilde{\beta}(z; \mathbf{x})^3 \right) \quad . \quad (125)$$

⌘ In Eq.(124), there could appear the single remaining connected two-loop vacuum diagram , with its own symmetry factor $1/8$. Since it is one-particle irreducible, it vanishes for diaphonies; not, however, for the χ^2 discrepancy that we shall discuss below (*cf* sect.6.4).

6.3.7 The two-point function

As mentioned in sect. 5.1.2 we still have to find the two-point correlation $F_2(t; \mathbf{x}_1, \mathbf{x}_2)$. This can also be done diagrammatically, by computing the averages of T^k keeping \mathbf{x}_1 and \mathbf{x}_2 fixed and integrating over the $N - 2$ other points. The only nonvanishing extra diagrams are of the form

$$\times \text{---} \bullet \text{---} \bullet \text{---} \times \quad \rightarrow \quad \times \text{-----} \times = \tilde{\beta}(z; \mathbf{x}_1 - \mathbf{x}_2) \quad , \quad (126)$$

that again carry an effective factor $1/N$ because these diagrams also have one vertex less than they have bare propagators. We can therefore write

$$\langle e^{zT} \rangle_{\mathbf{x}_{1,2} \text{ fixed}} = \Omega(z) \left(1 + \frac{1}{N} \tilde{\beta}(z; \mathbf{x}_1 - \mathbf{x}_2) \right) \quad (127)$$

and then

$$F_2(t; \mathbf{x}_1, \mathbf{x}_2) = - \left(\int dz e^{-tz} \Omega(z) \tilde{\beta}(z; \mathbf{x}_1, \mathbf{x}_2) \right) \left(\int dz e^{-tz} \Omega(z) \right)^{-1} . \quad (128)$$

In the saddle-point approximation this become quite simple:

$$F_2(t; \mathbf{x}_1, \mathbf{x}_2) = -\tilde{\beta}(z_0(t); \mathbf{x}_1, \mathbf{x}_2) . \quad (129)$$

For the one-dimensional Euler diaphony we find

$$\tilde{\beta}_E(z; x) = \sum_{n \neq 0} \frac{6z}{\pi^2 n^2 - 6z} \exp(2i\pi n x) \quad (130)$$

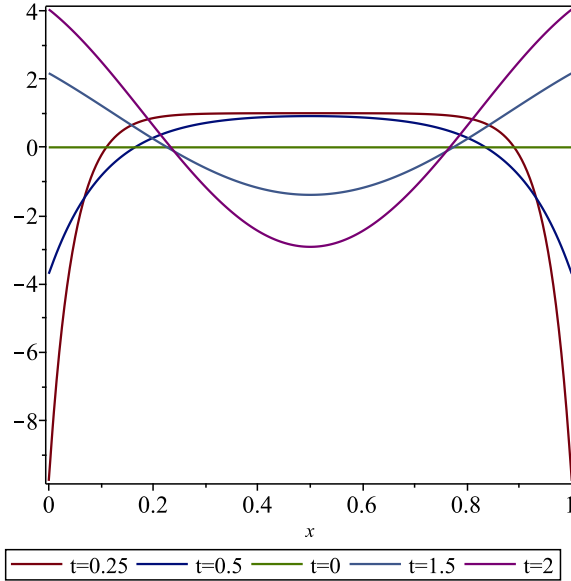
which satisfies the differential equation for $0 < x < 1$:

$$\frac{\partial^2}{(\partial x)^2} \tilde{\beta}_E(z; x) + 24z^2 \tilde{\beta} - E(z; x) = 24z^2 . \quad (131)$$

It has the solutions⁶⁰

$$\tilde{\beta}_E(z; x) = \begin{cases} 1 - \frac{y}{2} (e^{-yx} + e^{-y(1-x)}) / (1 - e^{-y}) & , \quad z < 0 \\ 1 - \frac{y}{2} (\sin(yx) + \sin(y(1-x))) / (1 - \cos(y)) & , \quad z > 0 \end{cases} , \quad (132)$$

where $y = \sqrt{24|z|}$.



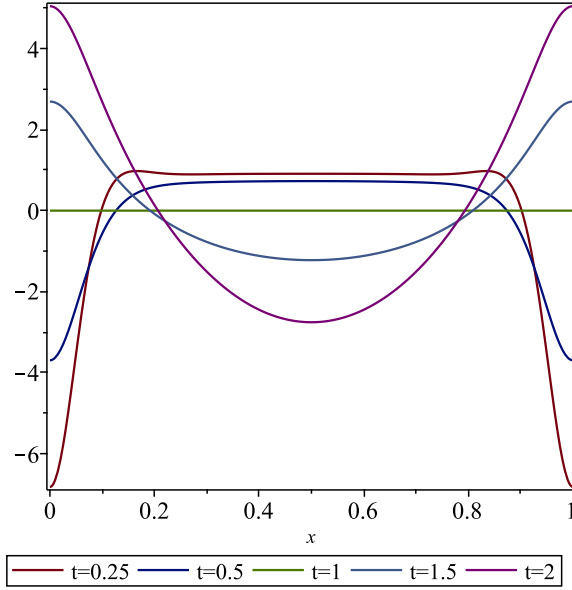
Here we plot the function $\tilde{\beta}_E(z_0(t); x)$ for various values of t . In the saddle-point approximation, $\tilde{\beta}(z(t); x_1 - x_2)$ is equal to $-F_2(t, x_1, x_2)$. For small t the ‘repulsion effect’ is evident. For large t , in contrast, there is ‘attraction’ and the points tend to cluster.

For the Block diaphony, where we have the saddle-point approximation $t =$

⁶⁰For the boundary conditions, see appendix [11.0.7](#).

$p/(p - z(t))$,

$$\tilde{\beta}_B(z_0(t); \mathbf{x}) = (t - 1) \left(-1 + \prod_{\mu=1}^d \frac{\sin((2c + 1)\pi \mathbf{x}^\mu)}{\sin(\pi \mathbf{x}^\mu)} \right) . \quad (133)$$



Here we plot $\tilde{\beta}_G(z_0(t); x)$ in the saddle-point approximation for the one-dimensional Gulliver diaphony, with $s = 0.5$. This two-point function is not easily obtained in closed form: I have simply summed numerically to large enough n . Qualitatively it is quite similar to $\tilde{\beta}_E(z_0(t); x)$.

6.3.8 Testing too much: the Dirac limit

The various diaphonies can be considered test of equidistribution. As an illustrative example, the block diaphony T_B tests how well Fourier modes with $|n^\mu| \leq c$ are integrated. Surely, if we include more and more modes, the test will become more stringent? We have

$$\sigma(T_B)^2 = \langle T_B^2 \rangle - \langle T_B \rangle^2 = 2 \sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^4 = \frac{2}{p} , \quad (134)$$

so that for c very large the variance of the T_B distribution vanishes: the distribution becomes a Dirac delta, and *every* point set⁶¹ ends up with $T_B \approx 1$. For the Gulliver and Euler diaphonies we find, similarly:

$$\sigma(T_G)^2 \approx \left(\frac{1-s}{4} \right)^d \text{ for } s \text{ approaching } 1 ,$$

⁶¹Because every point set finds itself, eventually, in the ensemble of random iid point sets.

$$\sigma(T_E)^2 \approx (21/40)^d \quad \text{for } d \text{ very large} . \quad (135)$$

The result holds generally: if $\sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^2 = 1$, then $\sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^4$ will tend to zero unless a *finite* number of strengths completely out-dominate all the other ones [20]. The ‘ultimate test’ is no test at all; the message is that one should not test ‘ad infinitum’, but *when* to stop testing is not clear. In this respect, as in others, Monte Carlo is an art rather than a prescription.

⋈ The limit of ‘large number of modes’ has its own Central Limit theorem: the sums $\sum_{\mathbf{n}}^* \sigma_{\mathbf{n}}^{2k}$ approach zero ever faster for increasing k , and therefore the density $P(t)$ takes on a Gaussian form [20].

6.4 Measures of nonuniformity: χ^2

6.4.1 The χ^2 as a discrepancy

The well-known χ^2 density can also be cast in the language of problem classes. In this case we divide the I^d hypercube into B non-overlapping, but not necessarily simply connected or even connected, regions (‘bins’). We define

$$\theta_n(\mathbf{x}) = \begin{cases} 1 & , \quad \mathbf{x} \text{ inside bin } n \\ 0 & , \quad \mathbf{x} \text{ outside bin } n \end{cases} . \quad (136)$$

We have

$$\theta_n(\mathbf{x}) \theta_m(\mathbf{x}) = \delta_{mn} \theta_n(\mathbf{x}) \quad , \quad \sum_{n=1}^B \theta_n(\mathbf{x}) = 1 . \quad (137)$$

The volume of the bins is given by

$$v_n = \int d\mathbf{x} \theta_n(\mathbf{x}) \quad , \quad \sum_{n=1}^B v_n = 1 . \quad (138)$$

The ‘Lego’ problem class⁶² now consists of functions that are piecewise constant over the bins:

$$f(\mathbf{x}) = \sum_n \alpha_n \theta_n(\mathbf{x}) . \quad (139)$$

The ensemble measure is again Gaussian, with

$$\langle \alpha_n \rangle_L = 0 \quad , \quad \langle \alpha_m \alpha_n \rangle_L = \delta_{mn} \frac{1}{v_n} . \quad (140)$$

⁶²Because the functions look like the ‘Lego plots’ common in experimental analysis.

It is this choice that singles out the χ^2 distribution out of all possible ‘Lego’-like discrepancies. The integration error,

$$\eta = \frac{1}{N} \sum_{j=1}^N f(\mathbf{x}_j) - \int d\mathbf{x} f(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^B \alpha_n \sum_{j=1}^N (\theta_n(\mathbf{x}_j) - v_n) \quad , \quad (141)$$

has expected square

$$\langle \eta^2 \rangle_L = \frac{1}{N^2} \sum_{n=1}^B \frac{1}{v_n} \left(\sum_j (\theta_n(\mathbf{x}_j) - v_n) \right)^2 \quad , \quad (142)$$

and this leads us to propose the measure of nonuniformity to be

$$T_L(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^B \frac{1}{v_n} \left(\sum_j (\theta_n(\mathbf{x}_j) - v_n) \right)^2 \quad . \quad (143)$$

This is exactly the χ^2 of the point set \mathbf{X} tested against the hypothesis of uniform distribution over I^d .

⋈ The fact that the bins may have any shape, or consist of disconnected parts, is the reason for considering (almost) any empirical test of a PRNG as a χ^2 test of uniformity in a possibly many-dimensional space with weird-looking bins (*e.g.* the poker test).

6.4.2 Large- N results for χ^2

The bare propagator for T_L is

$$\beta_L(\mathbf{x}, \mathbf{y}) = -1 + \sum_n \theta_n(\mathbf{x}) \theta_n(\mathbf{y}) / v_n \quad . \quad (144)$$

It is again tadpole-free, $\int \beta(\mathbf{x}, \mathbf{y}) d\mathbf{y} = 0$; but it is not translation invariant and hence not one-vertex reducible except when all volumes are equal, $v_n = 1/M$. On the other hand, we have the nice property

$$\int d\mathbf{z} \beta_L(\mathbf{x}, \mathbf{z}) \beta_L(\mathbf{z}, \mathbf{y}) = \beta_L(\mathbf{x}, \mathbf{y}) \quad . \quad (145)$$

This means that all bracelets come out the same:

$$B_k = B_1 = \int d\mathbf{x} \beta(\mathbf{x}, \mathbf{x}) = M - 1 \quad . \quad (146)$$

We find immediately that, for $N \rightarrow \infty$ ⁶³,

$$\begin{aligned}\Omega_L(z) &= (1 - 2z)^{-(B-1)/2} , \\ P_L(t) &= \frac{1}{2\Gamma\left(\frac{B-1}{2}\right)} \left(\frac{t}{2}\right)^{(B-3)/2} \exp\left(-\frac{t}{2}\right) .\end{aligned}\quad (147)$$

The saddle-point is reached for

$$z_0(t) = \frac{1}{2} \left(1 - \frac{B-1}{t}\right) . \quad (148)$$

The number $B - 1$ is called the *number of degrees of freedom*. Note that for $B = 1$ we have only the single bin I^d itself with trivially $T_L(\mathbf{X}) = 0$ for all point sets, hence no degrees of freedom.

⌘ The χ^2 discrepancy has expectation value $B - 1$ rather than 1. If we renormalize by scaling t to $t/(B - 1)$ the probability density is exactly that for the Block diaphony with $p = (B - 1)/2$; a curious result since the two notions of nonuniformity are defined in totally different ways!

6.4.3 Two-point function and $1/N$ corrections for χ^2

Because of property (145) we can immediately find the dressed two-point function:

$$\begin{aligned}\times \text{-----} \times &= \tilde{\beta}_L(z; \mathbf{x}, \mathbf{y}) = \frac{2z}{1 - 2z} \beta_L(\mathbf{x}, \mathbf{y}) \\ &= \frac{2z}{1 - 2z} \left(-1 + \sum_n \frac{\theta_n(x)\theta_n(y)}{v_n} \right) .\end{aligned}\quad (149)$$

As before, we find repulsion for small t ⁶⁴. The $1/N$ correction terms now include an extra diagram: we have

$$\Omega_L(z) = (1 - 2z)^{-(M-1)/2} \left(1 + \frac{\mathcal{A}}{N}\right) ,$$

⁶³We take B to be odd for simplicity here.

⁶⁴In the saddle-point approximation.

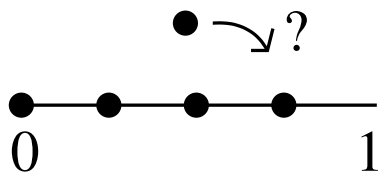
$$\begin{aligned}
\mathcal{A} &= -\frac{1}{8} \text{ (circle with one dot) }^2 - \frac{1}{4} \text{ (circle with two dots) } + \frac{1}{8} \text{ (two circles connected by a dot) } \\
&\quad + \frac{1}{12} \text{ (two circles connected by two dots) } + \frac{1}{8} \text{ (two circles connected by a dot and a line) } \\
&= \frac{z^2}{2(1-2z)^2} \left(\sum_n \frac{1}{v_n} - M^2 - 2M + 2 \right) \\
&\quad + \frac{z^3}{3(1-2z)^3} \left(5 \sum_n \frac{1}{v_n} - 3M^2 - 6M + 4 \right) . \quad (150)
\end{aligned}$$

7 Superuniform point sets

Point sets with a discrepancy/diaphony (considerably) lower than that expected for truly random ones are called superuniform.

7.1 Fixed point sets *vs* streams

We shall discuss several approaches to the construction of point sets with low discrepancy/diaphony. Here it becomes important to distinguish between fixed-size point sets and streams. If a point set of n points has very low diaphony, then adding an $(n+1)^{\text{th}}$ one will be problematic from the uniformity point of view. This point set of 4 points

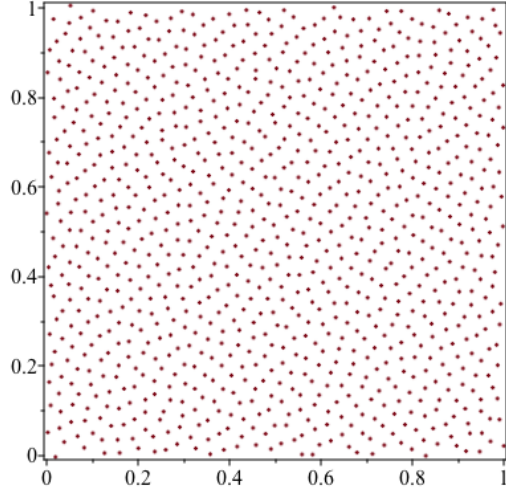


at $0, 1/4, 1/2, 3/4$ has the smallest diaphony possible. Where can we put the 5^{th} one without increasing the diaphony (note that points at 0 and at 1 coincide)? A point set of 5 points with minimal diaphony has its points at $0, 1/5, 2/5, 3/5, 4/5$.

7.1.1 Diaphony minimisation

For given N there exists the point set with minimal discrepancy/diaphony. *Finding* this point set is (in more than one dimension) not feasible. In general the best we can hope for is to obtain point sets with very low diaphony. Two strategies can be envisaged: either minimising the diaphony by shifting points around, or invoking some rule. Numerical minimisation typically relies on (a) descending methods that use expressions for the *gradient* of the diaphony⁶⁵, or (2) the Metropolis algorithm (*cf* sect.9.4.1). Both methods are very slow.

⁶⁵This cannot be done for discrepancy since the local discrepancy is by construction not differentiable where it counts.



The plot shows a two-dimensional low-discrepancy point set with $N = 1000$. It was obtained by student T.Blank in about a week's computing time, by descending the Gulliver diaphony with $s = 0.5$. The obtained diaphony is $7.2 \cdot 10^{-8}$. It is clear that simply running a PRNG many times and selecting the 'best' point set is a hopeless strategy.

7.1.2 Korobov sequences: good lattice points

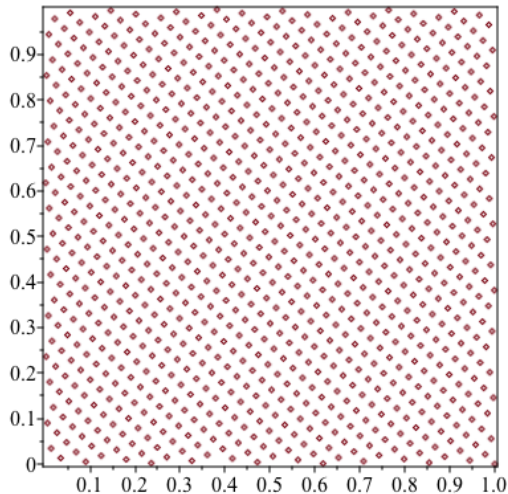
A widely used strategy is that of Korobov sequences, or the method of *good lattice points*. For a d -dimensional N -point set this consists of identifying a 'good' lattice vector with natural coefficients

$$\vec{g} = (g^1, g^2, \dots, g^d) , \quad (151)$$

and then the points are defined by

$$\mathbf{x}_k = \left(\left\{ \frac{k g^1}{N} \right\}, \left\{ \frac{k g^2}{N} \right\}, \dots, \left\{ \frac{k g^d}{N} \right\} \right) , \quad k = 1, 2, \dots, N . \quad (152)$$

We can take $g^1 = 1$ without loss of generality. The other components g^j should at least be relatively prime to N and to one another. A possibility



particular definition of nonuniformity. In my opinion fixed- N point sets are of limited use since for a realistic nontrivial integration problem it is not known *a priori* what N ought to be. An important insight, however, is the following: in the complete point set, the set of the j^{th} coordinates of the points, $(x_1^j, x_2^j, \dots, x_N^j)$, are precisely the minimal-diaphony, equidistant point sets in one dimension, thus explaining in some qualitative way the low diaphony of the full d -dimensional point set⁶⁶.

for $d = 2$ is this: if N equals the n^{th} Fibonacci number F_n we can take $\vec{g} = (1, F_{n-1})$. The plot shows the result for $F_{16} = 987$, $F_{15} = 610$. The uniformity is obvious, but so is the lack of rotational symmetry, especially when compared to the plot in sect.7.1.1. Essentially, such point sets have the same advantages and drawbacks as ‘hypercubic’ lattices. The low value of diaphony is, in some way, a result of specialising the lattice vector to *that*

7.2 QRNG algorithms

In view of the above, the more attractive idea is to search for low-diaphony *stream* algorithms, or QRNG: *Quasi-random number generators*.

7.2.1 Richtmeyer-Kronecker streams

Korobov sequences eventually ‘run out of steam’ since if we continue the rule (152) beyond $k = N$ the points will start to repeat. We can therefore envisage to let N approach infinity, and of course then the g^j have to approach infinity as well: the rational numbers g^j/N have to become *irrational*. This is the idea of Richtmeyer sequences: rather than identifying a lattice vector \vec{g} we

⁶⁶However, the vector $\vec{g} = (1, 1, \dots, 1)$ would give the same projections but an unacceptable more-dimensional set. Hence the requirement that the components of \vec{g} be mutually prime.

search for ‘irrational vectors’

$$\vec{g} = (\theta^1, \theta^2, \dots, \theta^d) , \quad (153)$$

where the numbers θ^j are all irrational numbers that are also mutually irrational⁶⁷. The quasi-random numbers \mathbf{x}_k are then given by

$$\mathbf{x}_k = (, \{k \theta^1\}, \{k \theta^2\} \dots, \{k \theta^d\}) . \quad (154)$$

The suitability of the irrationals θ^j can be investigated using their continued-fraction representation, that we shall now discuss.

7.2.2 Excursion into fractions (cont’d)

Let θ be a number in $(0, 1)$. Then

$$\frac{1}{\theta} = a + \theta' \quad , \quad a = \left\lfloor \frac{1}{\theta} \right\rfloor \quad (155)$$

so that a is an natural number and θ' is a number in $[0, 1)$. We can therefore repeat this procedure to arrive at the *continued-fraction representation* of θ :

$$\theta = \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \frac{1}{a_4 + \dots}}}} \equiv [a_1, a_2, a_3, a_4, \dots] . \quad (156)$$

If θ is a rational number, a_j will become infinite for some j and the fraction stops there; we can then write $\theta = [a_1, a_2, \dots, a_{j-1}]$. For irrational θ the continued-fraction representation continues forever. Some examples:

$$\begin{aligned} \pi - 3 &= [7, 15, 1, 292, 1, 1, 1, 2, 1, 3, 1, 14, 2, 1, 1, 2, 2, 2, \dots] , \\ (\sqrt{5} - 1)/2 &= [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, \dots] , \\ \sqrt{26} - 5 &= [10, 10, 10, 10, 10, 10, 10, 10, 10, 10, 10, 10, 10, \dots] , \\ \sqrt{3} - 1 &= [1, 2, 1, 2, 1, 2, 1, 2, 1, 2, 1, 2, 1, 2, 1, 2, 1, \dots] , \\ \sin(1) &= [1, 5, 3, 4, 19, 2, 2, 2, 2, 7, 2, 2, 1, 136, 3, 20, 3, 1, 3, \dots] , \\ 2^{1/3} - 1 &= [3, 1, 5, 1, 1, 4, 1, 1, 8, 1, 14, 1, 10, 2, 1, 4, 12, 2, 3, \dots] , \\ e - 2 &= [1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, 1, 1, 10, 1, 1, 12, 1, 1, 14, \dots] . \end{aligned} \quad (157)$$

⁶⁷That is, every ratio θ^i/θ^j is also irrational.

If the continued-fraction expansion is periodic, then θ will be the solution of a quadratic equation with integer coefficients, a *quadratic irrational*. Therefore all other irrational numbers, like $2^{1/3}$, have continued-fraction expansions that are *aperiodic*.

✂ The fact that the continued-fraction coefficients for most irrational numbers do not form a periodic pattern might lead one to propose these as a source of ‘truly random’ integers. This is a bad idea, since (a) the irrational number would have to be known to many millions of digits, and (b) non-periodicity does not imply randomness (*cf* Eq.(157) for $e-2$), or even a *known* distribution of integers.

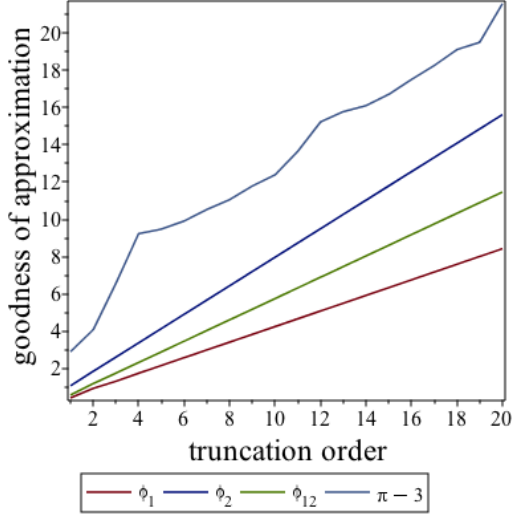
7.2.3 Rational approximations to irrationals

Truncating the continued-fraction representation gives us a method to approximate numbers by rational numbers:

$$\begin{aligned}\theta &\approx \theta_n = \frac{p_n}{q_n} \ , \\ p_0 &= 0 \ , \quad q_0 = 1 \ , \\ p_1 &= 1 \ , \quad q_1 = a_1 \ , \\ p_n &= a_n p_{n-1} + p_{n-2} \ , \quad q_n = a_n q_{n-1} + q_{n-2} \ .\end{aligned}\tag{158}$$

We see that if a_j becomes infinite, then θ is the rational number p_{j-1}/q_{j-1} . In some sense, therefore, the number π is *almost* rational since $a_4 = 292$ is so large. In fact, replacing 292 by infinity gives the Chinese approximation $\pi \sim 355/113$, which gets 6 decimals correct⁶⁸. We can also claim to know the *most irrational number in the universe*: it is that for which all the coefficients

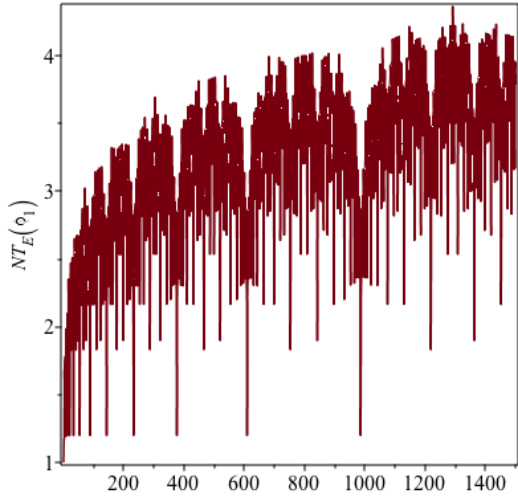
⁶⁸This fraction is called the *Milü*, established by Zu Chongzhi (429-500 AD).



a_j are 1, *i.e.* the golden ratio $\phi_1 \equiv (-1 + \sqrt{5})/2 = 0.618\dots$. This plot shows the goodness of the rational approximation (the number of correct decimal digits) for ϕ_1 , and also for $\phi_2 = [2, 2, 2, 2, \dots] = -1 + \sqrt{2}$, $\phi_{12} = [1, 2, 1, 2, 1, 2, \dots] = -1 + \sqrt{3}$, and $\pi - 3$. The almost-rationality of π is seen from the jump in goodness from $n = 3$ to $n = 4$. The smaller the continued-fraction coefficients, the worse the rational approximation: no curve exists below that for ϕ_1 .

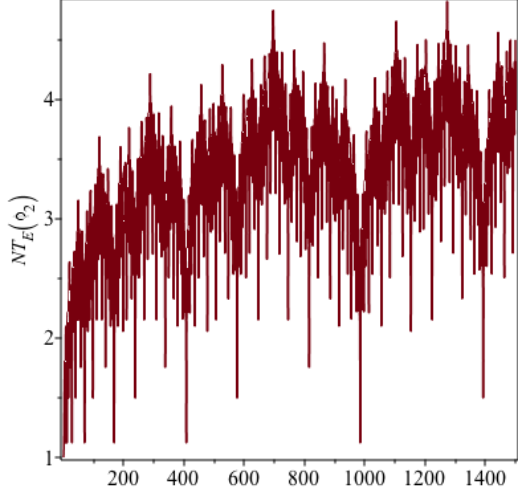
7.2.4 Almost-equidistancy for Richtmeyer sequences

Let us consider Richtmeyer sets x_k , $k = 1, \dots, N$, where $x_k = \{k\phi_1\}$, $\phi_1 = [1, 1, 1, 1, \dots]$ being the golden ratio. We plot the running value of $NT_E(\phi_1)$,

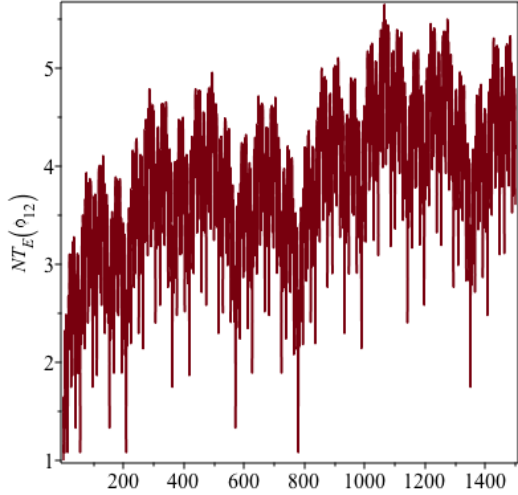


using the Euler diaphony T_E . The number ϕ_1 is approximated by ratios of Fibonacci numbers, F_{n-1}/F_n , the approximation improving for increasing n . Every time $N = F_m$ for some m , the distribution of points is *almost* equidistant with $x_k \sim q(k)/F_n$, where $q(k) \in [1, F_n]$ is some function of $k \in [1, F_n]$. That is, when $N = F_n$ the distribution of points has essentially the smallest possible diaphony.

This is the source of superuniformity: in between these ‘optimal’ values the diaphony cannot grow too much before coming down again. From appendix 11.0.8 we see that the Fibonacci numbers approximate $F_n \sim (1.618)^{-n}$ so that the diaphony moves out further and further as n increases. This is the reason



how close together the ‘optimal’ values are.



for the logarithmic term in the Roth bound. A similar phenomenon is observed for $NT_E(\phi_2)$ with $\phi_2 = [2, 2, 2, 2, \dots] = -1 + \sqrt{2}$, except that the returns to almost-equidistancy are now further apart, since $1/\phi_2 \sim 2.414$. As in the previous case a fractal pattern is evident, corresponding to a rational approximation to the irrational that is less than optimal. The quality of the diaphony is seen to depend on An intermediate case is that of a mixture of coefficients 1 and 2, for instance $\phi_{12} = [1, 2, 1, 2, \dots] = -1 + \sqrt{3}$ which is plotted here. The ‘optimal’ values (see appendix 11.0.8) are spaced $\sim (1.932)^n$, and the diaphony is minimal for odd n , next-to-minimal for even n . The minima are closer together than for ϕ_2 . We observe similar behaviour for other irrationals such as $[1, 1, 2, 1, 1, 2, 1, 1, 2, 1, 1, \dots] = -1 + \sqrt{5/2}$.

7.2.5 van der Corput streams

In the above Richtmeyer sequences, the (one-dimensional) distribution of points is *almost* equidistant at the ‘optimal’ values of N . We can improve on that using the following approach. The *van der Corput transform* $\phi_b(n)$ of an integer $n \geq 0$ in base b is defined as follows: if n has the b -ary expansion

$$n = n_0 + n_1b + n_2b^2 + n_3b^3 + \dots \quad (159)$$

then

$$\phi_b(n) = n_0 b^{-1} + n_1 b^{-2} + n_2 b^{-3} + n_4 b^{-4} + \dots \quad (160)$$

An explicit algorithm is given here: for n given, it runs in time $\mathcal{O}(\log(n))$.

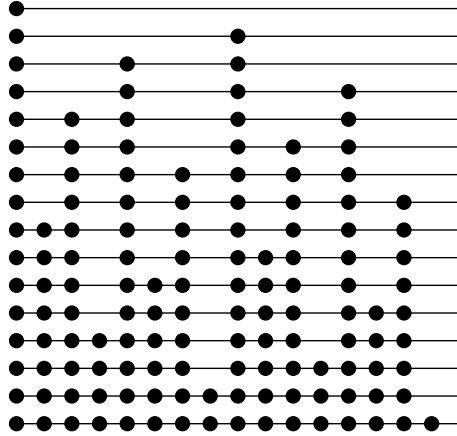
Algorithm 4 The van der Corput transform $\phi_b(n)$

```

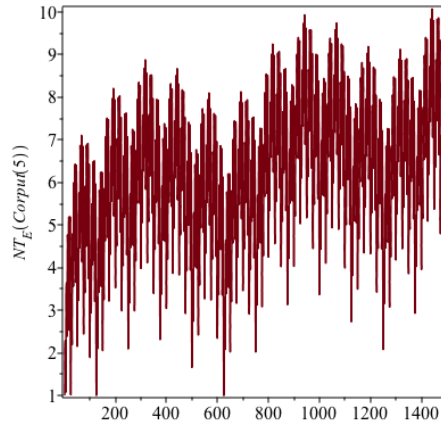
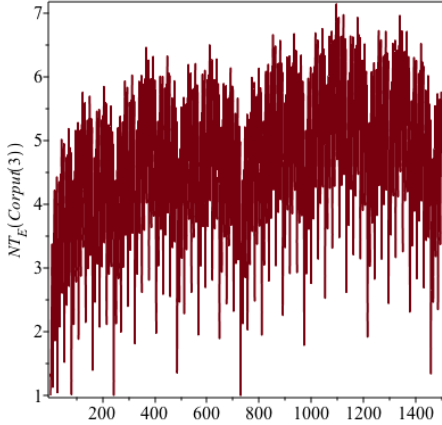
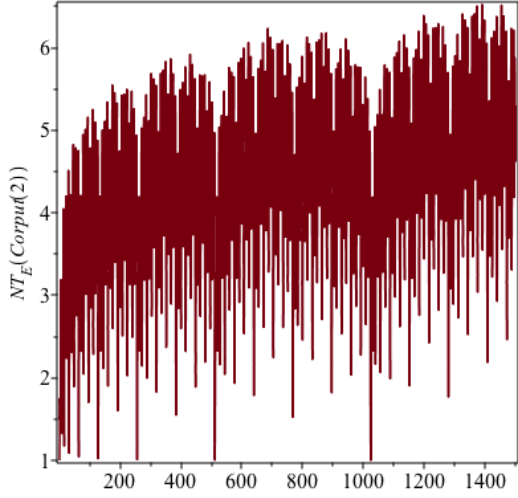
{The base is  $b$ , the input number is  $n > 0$ . }
 $m \leftarrow n$  {this avoids changes in  $n$ }
 $k \leftarrow 0$ 
while  $m > 0$  do
   $k \leftarrow k + 1$  ,  $c_k \leftarrow m \bmod b$  {get the next digit}
   $m \leftarrow (m - c_k)/b$ 
end while
 $j \leftarrow k$  ;  $\phi \leftarrow 0$  {start at the least significant digit;  $\phi$  is the output}
while  $j > 0$  do
   $j \leftarrow j - 1$  ,  $\phi \leftarrow (c_j + \phi)/b$ 
end while
return  $\phi$ 

```

In one dimension, the point set $\{\phi_b(1), \phi_b(2), \dots, \phi_b(N)\}$ will be exactly equidistant if $N = b^p$ for some power p : we return to the ‘best’ case as



often as possible for given b . Here we represent how the first van der Corput numbers $\phi_2(n)$ ($= 0, 2, \dots, 15$) are filling the unit interval $[0, 1)$. At line 1 ($n = 0$) the discrepancy is automatically minimal. In line 2 the 2 points are spaced equidistantly, as they are in line 4, 8, and 16. Every time 2^k points have been filled in the discrepancy returns to its minimal value. In between, at line



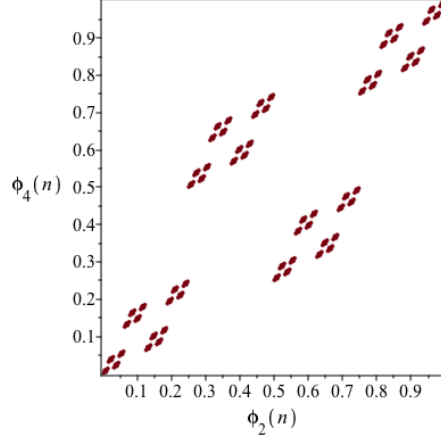
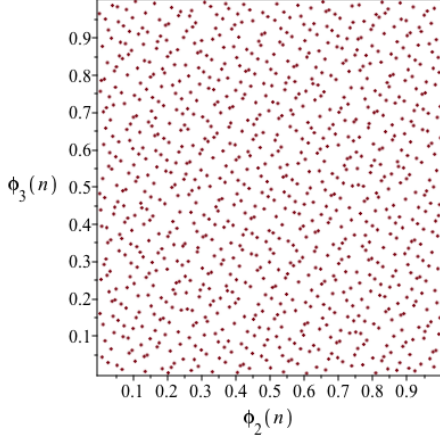
6, say, the discrepancy is ‘almost optimal’. This is evidenced by looking at $NT_E(\mathbf{X})$ as before. The evolution of the Euler diaphony shows an even more explicitly fractal pattern, and the optimal values at $N = 2^p$ have the absolute minimal diaphony. Inbetween, for $N = 2^p + 2^{p-1}$, say, the diaphony is slightly higher, and for $N = 2^p + 2^{p-2} \pm 2^{p-2}$ it is higher again. For larger values of the base b the minima are necessarily spread further apart. Below we give similar results for $b = 3$ and $b = 5$.

7.2.6 Van der Corput sequences in more dimensions

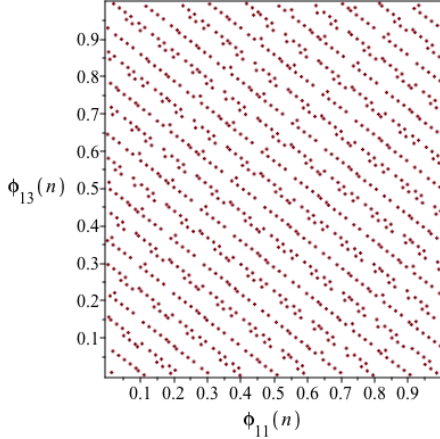
In two dimensions we may consider the sequences

$$\mathbf{x}_n = (\phi_2(n), \phi_3(n)) \quad \text{and} \quad \mathbf{x}_n = (\phi_2(n), \phi_4(n)) \quad . \quad (161)$$

The first 1000 points are displayed below.



Whereas the choice $b = 2, 3$ is quite acceptable for a superuniform point set, the choice $b = 2, 4$ is extremely unlucky: the pattern is fractal, with



diamonds built up from smaller diamonds and so on. It is clear that in more dimensions the bases must be relatively prime. But that implies that with increasing dimension d the bases must increase *at least as fast* as the lowest d primes. Going back to the one-dimensional projections of the point set, that implies that some ‘optimal’ values are going to be very far apart for appreciable d . As an example, we plot the point set for $b = 11, 13$, which is a two-dimensional

projection for all van der Corput sequences for $d \geq 6$. This shows that the discrepancy/diaphony can become large. In [21] the following result is given:

$$L_2^* < \frac{(\log N)^{2d}}{N^2} K_d \quad , \quad K_d = \prod_{j=1}^d \frac{(b_j - 1)^2}{\log(b_j)} \quad . \quad (162)$$

so we can determine the value of N_d for which the van der Corput sequence

d	K_d	N_d
2	5.253	$0.1560 \cdot 10^7$
3	52.24	$0.1807 \cdot 10^{12}$
4	966.4	$0.1079 \cdot 10^{18}$
5	$0.4030 \cdot 10^5$	$0.3192 \cdot 10^{24}$
6	$0.2263 \cdot 10^7$	$0.2144 \cdot 10^{31}$
7	$0.2045 \cdot 10^9$	$0.3649 \cdot 10^{38}$
8	$0.2252 \cdot 10^{11}$	$0.1087 \cdot 10^{46}$

is *guaranteed* to have a discrepancy smaller than the ‘random value’ $(2^{-d} - 3^{-d})/N$. This is given in the table. Since we only consider an *upper limit* on L_2^* here, the picture is not really so bleak. Nevertheless, the fact that *asymptotically* van der Corput sequences are superuniform does not guarantee useful behaviour for moderate N .

7.2.7 Niederreiter streams

If the van der Corput sequences deteriorate in higher dimensions because the bases b become large, we may decide to keep the same base b in all d dimensions but change the ordering of the points. This is the idea behind the *Niederreiter sequences*. For given base b (think of $b = 2$) we find functions $p_j(k)$, ($j = 1, 2, \dots, b$) with the following property: for all m , if k runs from 1 to b^m then $p_j(k)$ takes on all values between 1 and b^m as well. In other words, $p_j(k)$ is a permutation of the numbers $(1, 2, \dots, b^m)$. Having found these, the Niederreiter sequence is defined by

$$\mathbf{x}_k = \left(\phi_b(p_1(k)), \phi_b(p_2(k)), \dots, \phi_b(p_d(k)) \right) . \quad (163)$$

The trick is of course in finding the permutation functions p_j . A method of doing so is described in [22] and implemented independently in [23]. The (asymptotic!) bounds on the discrepancy are very much better than those for van der Corput sequences.

8 Variance reduction

8.1 Stratified sampling

8.1.1 General strategy

8.1.2 An example: VEGAS

8.1.3 An example: PARNI

8.2 Importance sampling

8.2.1 General strategy

8.2.2 Multichanneling

9 Non-uniform PRNGs

Often (and especially in the case of importance sampling) we are required to generate *non*-uniform pseudorandom numbers, and for many given densities a great number of strategies and tricks have been developed.

9.1 The Art of Transforming, Rejecting, and Being Smart

In generating nonuniform random numbers one invariably starts out with a source of iid pseudorandom numbers uniform in $(0, 1)$ ⁶⁹. These are then subjected to transformations, decisions and combinations. These can be broadly classified under

- Inversions (mappings) where a function of a variate is computed;
- Rejections where variates are kept or rejected according to some criterium;
- Cleverly combining several variates into new ones;
- The building of a repertoire of tricks as a basis for new tricks;
- Random-walk methods that employ ergodicity.

We shall discuss examples of all these. *The* reference text is the book by Devroye [24].

9.2 The UA formalism

9.2.1 Unitary algorithms as words and as pseudocode

Experience shows that, when we perform all kind of manipulations on random variates, it quickly becomes unclear what is precisely the density of the resulting numbers. Here it becomes useful to adhere to the *Unitary Algorithm* (UA) formalism. This is actually nothing more than integration statements with the number 1 on the left-hand side of the equation. These statements

⁶⁹Some care has to be taken here since quite often the special values 0 or 1 can give rise to numerical problems. We shall assume that they are never generated by your favourite PRNG.

can be put in words as well, and describe immediately (pseudo)code for software implementation of algorithms. The simplest example is

$$1 = \int_0^1 d\rho . \quad (164)$$

In words this reads ‘there exists an algorithm for generating numbers ρ uniformly between 0 and 1’. And this is true, since we assume that we have a good PRNG at hand. In pseudocode, the UA statement reads

$\rho \leftarrow \mathbf{prng}$

which just means ‘get a number out of your PRNG’. The second ingredient of the UA formalism is the multiplication by unity in the form of saturated Dirac deltas. For instance we may extend Eq.(164) as follows:

$$\begin{aligned} 1 &= \int_0^1 d\rho \int dx \delta(x - \rho^2) \\ &= \int_0^1 d\rho \int dx \frac{1}{2\rho} \delta(\rho - \sqrt{x}) \\ &= \int dx \frac{1}{2\sqrt{x}} \theta(0 < \sqrt{x} < 1) \\ &= \int_0^1 dx \frac{1}{2\sqrt{x}} . \end{aligned} \quad (165)$$

In words, we now have the statement ‘there is an algorithm for generating the density $1/(2x^{1/2})$ between 0 and 1’, and indeed there is: in pseudocode it reads

$\rho \leftarrow \mathbf{prng}$
 $x \leftarrow \rho^2$

In the UA description, the 1 on the left is to ensure that we are dealing with properly normalised densities. For instance, $\theta(0 < x < 1/2)$ is not a density, but $2\theta(0 < x < 1/2)$ is. Furthermore, in the step $\int dx \delta(x - \rho^2)$ the boundaries on x are $\pm\infty$ so that there *always* is some x for any ρ : this ensures that the algorithm actually finishes and yields a result. The final boundaries 0 and 1 on x arise from the fact that when we eliminate ρ from the Dirac delta we must make sure that ρ actually runs from 0 to 1.

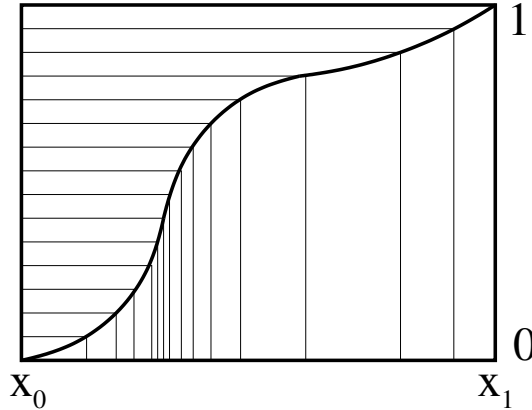
9.2.2 Inversion of variates in UA

One of the basic tools of the trade is the following. Suppose that we want to generate x according to some (properly normalised) density $g(x)$, between x_0 and x_1 . Suppose, furthermore, that we can find the primitive of $g(x)$:

$$G(x) = \int_{x_0}^x dy g(y) \quad \rightarrow \quad G(x_0) = 0; , G(x_1) = 1 \quad . \quad (166)$$

And suppose, *in addition*, that we are able to invert G somehow, so that $G^{-1}(z)$ can be computed. We then put x equal to $G^{-1}(\rho)$, in UA speak:

$$\begin{aligned} 1 &= \int_0^1 d\rho \int dx \delta(x - G^{-1}(\rho)) \\ &= \int dx \int_0^1 d\rho \delta(\rho - G(x)) \frac{1}{(G^{-1})'(\rho)} \\ &= \int dx \frac{1}{(G^{-1})'(G(x))} \theta(0 < G(x) < 1) \\ &= \int_{x_0}^{x_1} dx g(x) \end{aligned} \quad (167)$$



An illustration of the transformation method. We plot a function $G(x)$ and intersect it with regularly spaced horizontal lines. The resulting x values are on the lower axis. Obviously, the *density* of x values must be proportional to the *slope* of $G(x)$.

This method has the advantage that it uses only one **prng** per call, and it is very elegant. On the other hand, computing G^{-1} may not be possible analytically, or be very time-consuming; moreover it is really suited only to

one-dimensional densities so that for more-dimensional distributions one has to be able to successively integrate the density over its variables, and then also be able to invert the primitives. This can quickly become unfeasible. But for often-occurring, simple distributions this method works like a dream. Some examples:

```
 $\rho \leftarrow \text{prng}$   
 $x \leftarrow -\log(\rho)$ 
```

generates the exponential density $\exp(-x)\theta(x > 0)$, and

```
 $\rho \leftarrow \text{prng}$   
 $x \leftarrow \arctan(\pi(\rho - 1/2))$ 
```

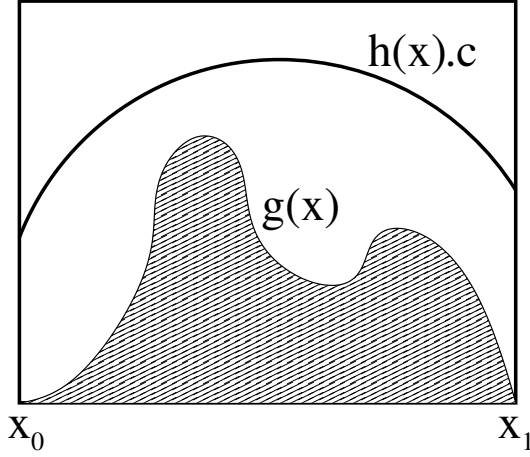
generates the Cauchy density $\pi^{-1}/(1 + x^2)$.

9.2.3 Rejection of variates in UA

This method works for more-dimensional densities as well as for one-dimensional ones. We want to generate a difficult target density $g(\mathbf{x})$ in some region Γ . ‘Difficult’ here means that we cannot use inversion, in fact we may not even know the normalisation of $g(\mathbf{x})$. Suppose that we *are* able to generate \mathbf{x} according to *another* density $h(\mathbf{x})$, and that $g(\mathbf{x}) < c h(\mathbf{x})$ in Γ for some known number⁷⁰ c . We then use the following algorithm, in pseudocode:

```
loop  
  generate  $\mathbf{x}$  according to  $h(\mathbf{x})$   
   $\rho \leftarrow \text{prng}$   
  if  $\rho < g(\mathbf{x})/(c h(\mathbf{x}))$  then  
    return {The value  $\mathbf{x}$  is accepted}  
  end if  
end loop
```

⁷⁰Larger than 1 if both g and h are properly normalised.



This illustrates the rejection method in one dimension. We want to fill the histogram of $g(\mathbf{x})$, the shaded area, uniformly. The algorithm relies on filling the area below $ch(\mathbf{x})$ uniformly, and cutting away, on a probabilistic basis, the points above the shaded area. That is, points are accepted with a probability $g(\mathbf{x})/(ch(\mathbf{x}))$.

We can analyse this algorithm the UA way as well, in spite of the fact that we may have to go through the loop an unbounded number of times.

$$\begin{aligned}
1 &= \int_{\Gamma} d\mathbf{x} P(\mathbf{x}) , \\
P(\mathbf{x}) &= \int_{\Gamma} d\mathbf{y}_1 h(\mathbf{y}_1) \int_0^1 d\rho_1 \left[\theta \left(\rho_1 < \frac{g(\mathbf{y}_1)}{ch(\mathbf{y}_1)} \right) \delta(\mathbf{x} - \mathbf{y}_1) \right. \\
&\quad + \theta \left(\rho_1 > \frac{g(\mathbf{y}_1)}{ch(\mathbf{y}_1)} \right) \int_{\Gamma} d\mathbf{y}_2 h(\mathbf{y}_2) \int_0^1 d\rho_2 \left[\right. \\
&\quad \left. \theta \left(\rho_2 < \frac{g(\mathbf{y}_2)}{ch(\mathbf{y}_2)} \right) \delta(\mathbf{x} - \mathbf{y}_2) + \right. \\
&\quad \left. \left. \theta \left(\rho_2 > \frac{g(\mathbf{y}_2)}{ch(\mathbf{y}_2)} \right) \int_{\Gamma} d\mathbf{y}_3 h(\mathbf{y}_3) \int_0^1 d\rho_3 \cdots \right] \right] \quad (168)
\end{aligned}$$

The expression for $P(\mathbf{x})$ is potentially infinite, but we can *telescope* it⁷¹:

$$\begin{aligned}
P(\mathbf{x}) &= \int_{\Gamma} d\mathbf{y}_1 h(\mathbf{y}_1) \int_0^1 d\rho_1 \left[\theta \left(\rho_1 < \frac{g(\mathbf{y}_1)}{ch(\mathbf{y}_1)} \right) \delta(\mathbf{x} - \mathbf{y}_1) \right. \\
&\quad \left. + \theta \left(\rho_1 > \frac{g(\mathbf{y}_1)}{ch(\mathbf{y}_1)} \right) P(\mathbf{x}) \right]
\end{aligned}$$

⁷¹Remember the continued fractions?

$$\begin{aligned}
&= \int_{\Gamma} d\mathbf{y}_1 h(\mathbf{y}_1) \left[\frac{g(\mathbf{y}_1)}{c h(\mathbf{y}_1)} \delta(\mathbf{x} - \mathbf{y}_1) + \left(1 - \frac{g(\mathbf{y}_1)}{c h(\mathbf{y}_1)} \right) P(\mathbf{x}) \right] \\
&= \frac{g(\mathbf{x})}{c} + P(\mathbf{x}) - P(\mathbf{x}) \int_{\Gamma} d\mathbf{y}_1 \frac{g(\mathbf{y}_1)}{c} .
\end{aligned} \tag{169}$$

And we finally arrive at

$$P(\mathbf{x}) = g(\mathbf{x}) \left(\int_{\Gamma} d\mathbf{y} g(\mathbf{y}) \right)^{-1} , \tag{170}$$

so that the rejection algorithm is *self-normalising*. Also note that it is really necessary that $g(\mathbf{x}) < c h(\mathbf{x})$ everywhere. The rejection algorithm is conceptually very simple but also has its drawbacks. It uses at least two calls to **prng** to obtain a single variate. It is common to use a uniform distribution for $h(\mathbf{x})$ but that will not work if Γ is infinitely large. Finding a good value for c may not be easy, or impossible if $g(\mathbf{x})$ goes to infinity somewhere. Of course the algorithm will work for any sufficiently large c but its efficiency will become low if c is very large. Rejection is typically used if we have *almost* the distribution we want, and has only to be massaged a bit. The canonical reference to the method is to von Neumann [25], but surely the method must be older.

9.3 Repertoire and the Rule of Nifty

An important rôle in the generation of nonuniform variates is played by the buildup of a repertoire of algorithms, and by application of the following Rule of Nifty: *a clever way of constructively computing the normalisation of a density will usually lead to an efficient algorithm for generating the density.*

9.3.1 Building up a repertoire

It pays to keep in mind tricks that work. We give an example here. Suppose that we have an algorithm for generating $g(x)\theta(x > 0)$. Then, multiplying with a **prng** results in

$$1 = \int_0^{\infty} dy g(y) \int_0^1 d\rho \int dx \delta(x - \rho y)$$

$$\begin{aligned}
&= \int dx \int_0^\infty dy \frac{g(y)}{y} \int_0^1 d\rho \delta\left(\rho - \frac{x}{y}\right) \\
&= \int dx \int_0^\infty dy \frac{g(y)}{y} \theta(x < y) \\
&= \int_0^\infty dx P(x) \quad , \quad P(x) = \int_x^\infty dy g(y)/y \quad . \tag{171}
\end{aligned}$$

By repeated applications, we then establish that

$$\begin{aligned}
\rho_{1,2,\dots,k+1} &\leftarrow \mathbf{prng} \\
x &\leftarrow \rho_1 \rho_2 \cdots \rho_k \rho_{k+1}
\end{aligned}$$

gives the density $\log(1/x)^k/k!$ between 0 and 1. Note that this is not possible with rejection, and unfeasible with inversion. Similarly,

$$\begin{aligned}
\rho_{1,2,\dots,k+1} &\leftarrow \mathbf{prng} \\
x &\leftarrow -\log(\rho_1 \rho_2 \cdots \rho_k \rho_{k+1})
\end{aligned}$$

results in the density $x^k \exp(-x)/k!$. Weird densities are possible:

$$\begin{aligned}
\rho_{1,2} &\leftarrow \mathbf{prng} \\
x &\leftarrow -\rho_1 \log(\rho_2)
\end{aligned}$$

is what you need if you ever have to generate the *exponential integral* $E_1(x)$ [\[26\]](#); another weirdo algorithm is

$$\begin{aligned}
\rho_{1,2} &\leftarrow \mathbf{prng} \\
x &\leftarrow \log(\rho_1) \log(\rho_2)
\end{aligned}$$

very useful if the need to generate $2K_0(2\sqrt{x})$ ever crosses your path⁷².

9.3.2 The normal distribution: the Box-Müller algorithm

As an example of the Rule of Nifty we can consider the Guassian, or normal density:

$$N(x) = \frac{1}{K} \exp(-x^2) \quad , \quad K = \int_{-\infty}^{\infty} dx \exp(-x^2) \quad . \tag{172}$$

Computing K by integrating the density in a straightforward way involves the error function, so generating normal variates by transformation calls for the *inverse error function*, which is horrible. Straightforward rejection from

⁷²The function K_0 is the modified Bessel function of the second kind of order zero [\[26\]](#).

a uniform distribution is also impossible on the interval $(-\infty, \infty)$. However, there exists the ‘doubling trick’, where polar coordinates are used:

$$\begin{aligned} K^2 &= \int dx dy \exp(-x^2 - y^2) \\ &= \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dr r \exp(-r^2) = \pi \int_0^{\infty} ds \exp(-s) = \pi . \end{aligned} \quad (173)$$

We already know how to generate the density $\exp(-s)$ (see sect.9.2.2). This method is known as the Box-Müller algorithm [27]. There are faster methods, but (to my mind) none so elegant.

Algorithm 5 The Box-Müller algorithm

```

 $\rho_{1,2} \leftarrow \text{prng}$ 
 $r^2 \leftarrow -\log(\rho_1)$ 
 $\phi \leftarrow 2\pi\rho_2$ 
 $x \leftarrow r \cos(\phi)$ 
 $y \leftarrow r \sin(\phi)$  { $x$  and  $y$  are independent normal variates}

```

9.3.3 The Euler algorithm

Quite often we are asked to generate numbers satisfying some constraint. An example of the Rule of Nifty is the *generalised Euler distribution*, the n -dimensional probability $P(x_1, x_2, \dots, x_n)$ given by

$$P(x_1, x_2, \dots, x_n) \sim x_1^{p_1} x_2^{p_2} \dots x_n^{p_n} \delta \left(\sum_{j=1}^n x_j - 1 \right) , \quad (174)$$

where $p_j \geq 0$ are integers. The normalisation factor is here unknown: we must also compute it. It might be tempting to generate the x_j in order, with x_j between 0 and $1 - x_1 - \dots - x_{j-1}$, but we can see straightaway that, first, the distribution of $x_{1,2,\dots,n-1}$ is independent of n (which is surely wrong) and, second, the value of x_n is completely fixed by $x_{1,2,\dots,n-1}$ so that its density cannot play a rôle. The better algorithm is given here, further on we shall discuss its various aspects.

Algorithm 6 The Euler density with parameters p_1, p_2, \dots, p_n

generate y_j in $(0, \infty)$ according to $\sim y^{p_j} \exp(-y)$ for $j = 1, 2, \dots, n$
 $s \leftarrow y_1 + y_2 + \dots + y_n$
 $x_j \leftarrow y_j/s$ for $j = 1, 2, \dots, n$

Let us write this out in the UA formalism. It reads

$$\begin{aligned}
1 &= \frac{1}{p_1! p_2! \dots p_n!} \int_0^\infty dy_1 \dots dy_n y_1^{p_1} e^{-y_1} \dots y_n^{p_n} e^{-y_n} \\
&\quad \int ds \delta(y_1 + \dots + y_n - s) \\
&\quad \int dx_1 \dots dx_n \delta\left(x_1 - \frac{y_1}{s}\right) \dots \delta\left(x_n - \frac{y_n}{s}\right) . \quad (175)
\end{aligned}$$

First, we eliminate the y 's in favor of the x 's, taking care to correctly handle the factors of s coming from the Dirac deltas, and then we do the integral over s :

$$\begin{aligned}
1 &= \frac{1}{p_1! p_2! \dots p_n!} \int_0^\infty dx_1 \dots dx_n x_1^{p_1} \dots x_n^{p_n} \exp(-s(x_1 + \dots + x_n)) \\
&\quad \int_0^\infty ds s^{p_1 + \dots + p_n + n} \delta(s(x_1 + \dots + x_n) - s) \\
&= \frac{1}{p_1! p_2! \dots p_n!} \int_0^\infty dx_1 \dots dx_n x_1^{p_1} \dots x_n^{p_n} \delta(x_1 + \dots + x_n - 1) \\
&\quad \int_0^\infty ds s^{p_1 + p_n + n - 1} e^{-s} \\
&= \frac{\Gamma(p_1 + \dots + p_n + n)}{\Gamma(p_1 + 1) \dots \Gamma(p_n + 1)} \int_0^\infty dx_1 \dots dx_n x_1^{p_1} \dots x_n^{p_n} \delta\left(\sum_j x_j - 1\right) . \quad (176)
\end{aligned}$$

We see that the density is precisely correct, and we obtain the normalisation into the bargain. The ‘original’ distribution of the y 's contains a factor $\exp(-y)$. Such a damping factor is necessary if we want to allow for arbitrarily large y values: these *must* be allowed since $(x_1, x_2, \dots, x_n) = (1, 0, \dots, 0)$ must be possible. The damping factor might also have been $\exp(-y^2)$, say, but our choice is seen to be the better one since it leads to a uniform sampling.

9.3.4 The Kinderman-Monahan algorithm

Let us generate two variates v, u iid uniformly, $0 < u < 1$ and $-1 < v < 1$, and consider their ratio:

$$\begin{aligned}
1 &= \frac{1}{2} \int_0^1 du \int_{-1}^1 dv \int dx \delta\left(x - \frac{v}{u}\right) \\
&= \frac{1}{2} \int dx \int_0^1 du u \int_{-1}^1 \delta(v - xu) \\
&= \frac{1}{4} \int dx \int_0^1 d(u^2) \theta(x^2 u^2 < 1) .
\end{aligned} \tag{177}$$

The step functions can be translated as

$$\begin{aligned}
&\theta(x^2 u^2 < 1) \theta(0 < u < 1) = \\
&\theta(|x| < 1) \theta(0 < u < 1) + \theta(|x| > 1) \theta(0 < u < 1/|x|) ,
\end{aligned} \tag{178}$$

and so we get

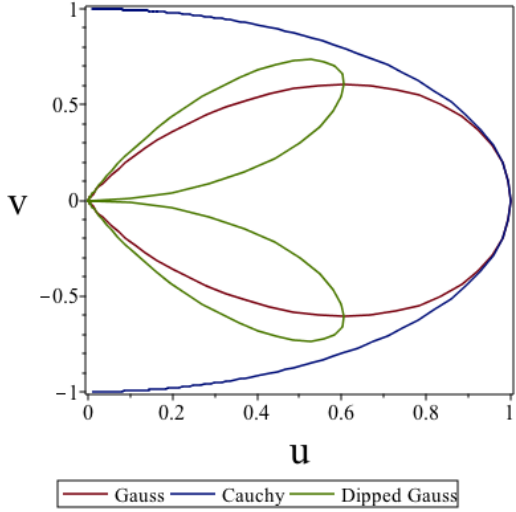
$$1 = \int dx h(x) \quad , \quad h(x) = \frac{1}{4} \left[\theta(|x| < 1) + \frac{1}{x^2} \theta(|x| > 1) \right] . \tag{179}$$

Thus we have a method to generate a density $h(x)$ that has a bump around zero and tails that fall off as $1/x^2$. We may use this one to generate other distributions by rejection⁷³: The Kinderman-Monahan, or *ratio of uniforms* algorithm[28] . Let us generate u and v and keep only points that fall inside a certain region defined by a function f :

$$\begin{aligned}
&\int_0^1 du \int_{-1}^1 dv \theta\left(u < \sqrt{f(v/u)}\right) \int dx \delta(x - v/u) \\
&= \int dx \int_0^1 du u \theta(u < \sqrt{f(x)}) \theta(-1 < ux < 1) \\
&\sim \int dx f(x) ,
\end{aligned} \tag{180}$$

⁷³Since rejection is self-normalising we can afford to be somewhat sloppy with the overall factors.

and this proves that x is generated proportional to $f(x)$. Note that we need $f(x) < 1$ as well as the fact that $\sqrt{f(x)} < 1/|x|$ in order for the last step function to be irrelevant. The Kinderman-Monahan algorithm essentially asks only for the determination that the points (u, v) are inside a certain region, and this can be made quite efficient.



The region of acceptance is bounded by a curve $(u(\tau), v(\tau))$ where $v(\tau) = \tau u(\tau)$, $u(\tau) = \sqrt{f(\tau)}$. Here we give three examples: the Gauss density $\sim \exp(-x^2)$, the Cauchy density $\sim 1/(1+x^2)$, and a ‘dipped Gauss’ density $\sim x^2 \exp(-x^2)$. We may relax the constraints $0 < u, |v| < 1$, and simply look for the maxima of $u(\tau)$ and $|v(\tau)|$ in order to decide on the rectangle in which to sample the points uniformly. For the Cauchy distribution, the

region of acceptance is a half circle, and this gives us an even simpler method of generating it, which avoids taking an arctangent at the cost of more random numbers.

Algorithm 7 Generating the Cauchy density by ratio of uniforms

```

loop
   $\rho_{1,2} \leftarrow \text{prng}$ 
   $r \leftarrow 2\rho_1 - 1$ 
  if  $r^2 + \rho_2^2 < 1$  then
    return  $r/\rho_2$ 
  end if
end loop

```

9.4 Random-walk algorithms

The algorithms discussed so far return iid variates. There is another strategy that is widely used, where new variates are determined from the previous ones. In the space of events the algorithm therefore performs a ‘random

walk' of which one must ensure that (a) the points visited have the desired density, also called the target density, (b) eventually the whole space is visited (that is, ergodicity holds), and (c) that the subsequent points can be argued to be *essentially* independent. Such algorithms also go under the name of Markov-Chain Monte Carlo (MCMC).

9.4.1 The Metropolis algorithm

This is nowadays also called the Metropolis-Hastings method [29]. The algorithm is actually quite simple. Suppose we are at a point \mathbf{x}_n . Then, using *some* prescription, we generate a *candidate* new point \mathbf{y} using a probability density $g(\mathbf{x}_n; \mathbf{y})$. We then compare the target densities $P(\mathbf{x})$ and $P(\mathbf{y})$:

$$R = \frac{P(\mathbf{y})}{P(\mathbf{x})} . \quad (181)$$

If $R > 1$ then we accept the candidate point $\mathbf{x}_{n+1} = \mathbf{y}$; this is reasonable since the candidate point is more probable than the old one. However, if $R < 1$ we *may* accept $\mathbf{x}_{n+1} = \mathbf{y}$, with probability R . If not, then we stick to the old point, and have $\mathbf{x}_{n+1} = \mathbf{x}_n$. The fact that we allow for a reduction in probability is what allows us to wander all over the space; otherwise we would simply be working our way towards a point of (locally) maximal probability. This method will work *provided* that we have *detailed balance*: we insist that

$$g(\mathbf{x}; \mathbf{y}) = g(\mathbf{y}; \mathbf{x}) . \quad (182)$$

The convergence of the Metropolis algorithm can best be pictured as follows. Suppose our space is populated by many random walkers, that at the start are distributed with some density $p_1(\mathbf{x})$. These walkers then each take one Metropolis step; some of them will remain where they are, others will be displaced. Their density after this first step is then $p_2(\mathbf{x})$, presumably a different one. This continues, so that as the algorithm proceeds we have a succession of densities $p_1(\mathbf{x}), p_2(\mathbf{x}), p_3(\mathbf{x}), \dots$. The idea is that $\lim_{n \rightarrow \infty} p_n(\mathbf{x}) = P(\mathbf{x})$. The UA formulation is

$$\begin{aligned} 1 &= \int d\mathbf{x} p_n(\mathbf{x}) , \\ p_n(\mathbf{x}) &= \int d\mathbf{z} p_{n-1}(\mathbf{z}) \int d\mathbf{y} g(\mathbf{z}; \mathbf{y}) \\ &\quad \times \left[\theta(P(\mathbf{y}) > P(\mathbf{z})) \delta(\mathbf{x} - \mathbf{y}) \right. \end{aligned}$$

$$\begin{aligned}
& + \theta(P(\mathbf{y}) < P(\mathbf{z})) \frac{P(\mathbf{y})}{P(\mathbf{z})} \delta(\mathbf{x} - \mathbf{y}) \\
& + \theta(P(\mathbf{y}) < P(\mathbf{z})) \left(1 - \frac{P(\mathbf{y})}{P(\mathbf{z})} \right) \delta(\mathbf{x} - \mathbf{z}) \Big] . \quad (183)
\end{aligned}$$

Let us now suppose that at some moment the walkers are actually distributed according to the target density: $p_{n-1}(\mathbf{x}) = P(\mathbf{x})$. After stepping, their density becomes

$$\begin{aligned}
p_n(\mathbf{x}) &= \int d\mathbf{z} P(\mathbf{z}) g(\mathbf{z}; \mathbf{x}) \theta(P(\mathbf{x}) > P(\mathbf{z})) \\
&+ \int d\mathbf{z} P(\mathbf{x}) g(\mathbf{z}; \mathbf{x}) \theta(P(\mathbf{x}) < P(\mathbf{z})) \\
&+ P(\mathbf{x}) \int d\mathbf{y} g(\mathbf{x}; \mathbf{y}) \theta(P(\mathbf{y}) < P(\mathbf{x})) \\
&- \int d\mathbf{y} g(\mathbf{x}; \mathbf{y}) P(\mathbf{y}) \theta(P(\mathbf{x}) > P(\mathbf{y})) . \quad (184)
\end{aligned}$$

Under detailed balance and renaming the integration variables, the first and fourth lines cancel and we are left with

$$\begin{aligned}
p_n(\mathbf{x}) &= P(\mathbf{x}) \int d\mathbf{z} g(\mathbf{x}; \mathbf{z}) \left[\theta(P(\mathbf{x}) < P(\mathbf{z})) + \theta(P(\mathbf{z}) < P(\mathbf{x})) \right] \\
&= P(\mathbf{x}) \int d\mathbf{z} g(\mathbf{x}; \mathbf{z}) = P(\mathbf{x}) . \quad (185)
\end{aligned}$$

We see that the target density $P(\mathbf{x})$ is a *fixed point* of the Metropolis algorithm, so that we can feel more or less confident that we shall approach our goal. Note the essential rôle played by the detailed balance requirement! For the rest we are completely free in the choice of g : we may change it at will at any moment. The choice of g *does* influence the performance of the algorithm, though. If we take ‘small’ steps the probability of accepting a candidate point is probably high, but it will take a (very) long time to cover the space. If the steps are ‘large’ then we move over the space quickly, but we may expect that not many of such proposed steps are accepted⁷⁴. A good rule of thumb seems to be that about half of the candidates should be accepted.

A drawback is the fact that the various points are *not* independent. This is typically overcome (hopefully) by taking a number of Metropolis steps

⁷⁴Especially if there are several probability maxima the chance of moving from the neighbourhood of one maximum to that of another can be very small.

before claiming the new point. How many steps are necessary and sufficient depends on the case: here, as everywhere, Monte Carlo is partly an art.

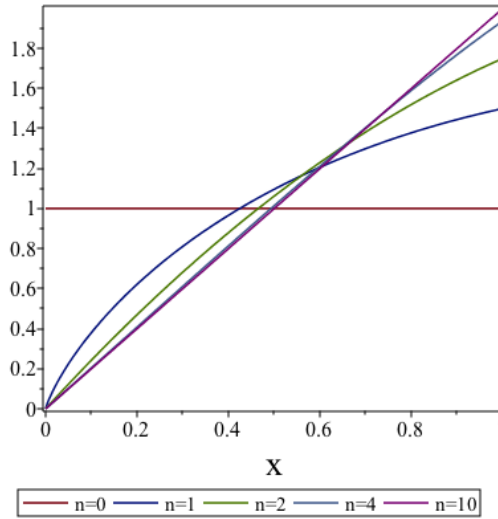
9.4.2 An elementary case study for Metropolis

We can completely analyse the Metropolis algorithm in a very simple case. This is the generation of the density $P(x) = 2x\theta(0 < x < 1)$. We shall use the step recipe $g(x; y) = \theta(0 < x, y < 1)$ which certainly gives detailed balance. The Metropolis step of Eq.(183) then takes the form

$$p_n(x) = \frac{x}{2} p_{n-1}(x) + \int_0^x dx p_{n-1}(z) + x \int_x^1 dz \frac{p_{n-1}(z)}{z} . \quad (186)$$

We start with $p_0(x) = \theta(0 < x < 1)$. The subsequent densities are

$$\begin{aligned} p_1(x) &= \frac{3}{2}x - x \log(x) , \\ p_n(x) &= \left(2 + \frac{1}{2^{n-1}(n-1)}\right)x - \frac{n+1}{2^n(n-1)}x^n , \quad n \geq 2 . \end{aligned} \quad (187)$$



We plot the shapes of $p_n(x)$ for $n = 0, 1, 2, 4$, and 10 . The density $p_{10}(x)$ is essentially indistinguishable from the target density $P(x) = 2x$. As we see from Eq.(187), the approach to P is exponentially fast. This way of generating $P(x)$ is of course much clumsier than the simple inversion rule $x \leftarrow \sqrt{\text{prng}}$ but it avoids taking a square root. Two-thirds of the candidates are accepted in this case.

An interesting observation is the following. Suppose that you would need to generate the distribution $x(3 - \log(x))/2$. Inversion seems pretty hopeless in this case, and rejection from a uniform density demands the computation of the logarithm. We see that an alternative method is to generate a random variable uniformly in $(0, 1)$ and then perform precisely *one* step of the Metropolis algorithm!

9.4.3 Applications of the Metropolis algorithm

The Metropolis algorithm finds a very natural (and, in fact, its original) application in statistical mechanics. For a system in the canonical ensemble its microconfigurations \mathbf{X} are distributed with a density

$$P(\mathbf{X}) \sim \exp\left(-\frac{1}{kT} H(\mathbf{X})\right) , \quad (188)$$

where k is Boltzmann's constant, T is the temperature, and $H(\mathbf{X})$ the Hamiltonian function of the system. Thermodynamic quantities are computed as averages over a sample of microconfigurations with this density. Calculating H is usually very cumbersome by itself, but in the Metropolis algorithm it is not necessary since what we are really interested in is the *ratio* of densities:

$$P(\mathbf{Y})/P(\mathbf{X}) = \exp\left(\frac{1}{kT} \left(H(\mathbf{X}) - H(\mathbf{Y})\right)\right) . \quad (189)$$

We therefore only have to compute the *change* in the Hamiltonian when stepping from \mathbf{X} to \mathbf{Y} . Especially for simple steps, the flipping of a single spin in an Ising system, say, this can be done very fast⁷⁵.

Another interesting application is that of *simulated annealing*, best exemplified by the Travelling Salesman problem: given a set of N ‘cities’ with known distances between them, find the shortest route that visits all the cities. For simplicity we demand that we come back to the starting city, and the direction of travel is unimportant. This gives us $(N - 1)!/2$ different routes, and a direct enumeration to find the optimal route becomes unpractical for $N > 13$ or so. What we can do instead is to assign a probability density to each route R :

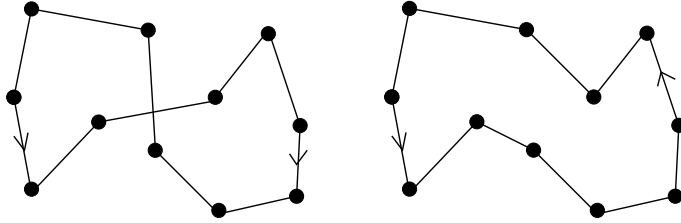
$$P(R) \sim \exp\left(-\frac{1}{kT} L(R)\right) , \quad (190)$$

where kT is as before, and $L(R)$ is the total length of route R . The ‘temperature’ is of course a totally fictional one. Using the Metropolis algorithm we

⁷⁵Flipping a single spin should not be the *only* possible step, however, especially close to an ordering transition. Occasionally, whole *blocks* of spins can be flipped as well. As we have already pointed out, choosing between different steps depending on the circumstances is perfectly allowed as long as detailed balance is maintained.

can sample this density; and by starting at high T and then gradually and carefully lowering it towards zero we may hope to end up in a maximum of the probability, *i.e.* a very short (or even *the* shortest) route. Typical steps in this game are the interchange of two cities, and (somewhat surprisingly at first sight) the reversal of a whole sequence of cities. Because of the triangle

inequality a route that crosses itself is always shortened by undoing the crossing, as indicated here. But that implies that the cities on one side of the crossing must be travelled in the *opposite* order!



9.4.4 Gibbs sampling

Another example of a random-walk method is Gibbs sampling, that aims at replacing the generation of a multidimensional density by a succession of one-dimensional ones. It is best described by an example. Suppose that we want to generate a three-dimensional density $P(\mathbf{x}) = P(x^1, x^2, x^3)$. This may not be possible directly, but suppose that we *can* generate x^1 for fixed $x^{2,3}$, and x^2 for fixed $x^{1,3}$, and also x^3 for $x^{1,2}$ fixed. Since these are one-dimensional densities, that is usually easier. Note that we have to normalize properly: for the generation of x^j the appropriate density is not $P(\mathbf{x})$ itself but P_j where (in the n -dimensional case)

$$P_j(x^1, \dots, x^{j-1}, x^j, x^{j+1}, \dots, x^n) = P(\mathbf{x}) \left[\int d\xi P(x^1, \dots, x^{j-1}, \xi, x^{j+1}, \dots, x^n) \right]^{-1}. \quad (191)$$

Now assume that at some step n in the Gibbs algorithm we have generated a point \mathbf{x} . A Gibbs step to the new point \mathbf{y} is then performed by first generating y^1 using $P(y^1, x^2, x^3)$, then generating y^2 according to $P_2(y^1, y^2, x^3)$, and finally generating y^3 from $P_3(y^1, y^2, y^3)$. At each step one new coordinate replaces an old one. Again considering a distribution of random walkers on our space, with density $p_n(\mathbf{x})$, we can then UA towards the density after one

Gibbs step:

$$\begin{aligned}
1 &= \int d^3\mathbf{x} p_n(\mathbf{x}) \\
&= \int dx^1 dx^2 dx^3 p_n(x^1, x^2, x^3) \\
&\quad \times dy^1 P_1(y^1, x^2, x^3) dy^2 P_2(y^1, y^2, x^3) dy^3 P_3(y^1, y^2, y^3) \\
&= \int d\mathbf{y} P(\mathbf{y}) R_n(\mathbf{y}) ,
\end{aligned} \tag{192}$$

where after reshuffling of factors we have

$$R_n(\mathbf{y}) = \int dx^1 \frac{p_n(x^1, x^2, x^3)}{\int d\xi P(\xi, x^2, x^3)} dx^2 \frac{P_1(y^1, x^2, x^3)}{\int d\xi P(y^1, \xi, x^3)} dx^3 \frac{P(y^1, x^2, x^3)}{\int d\xi P(y^1, y^2, \xi)} . \tag{193}$$

The new density $p_{n+1}(\mathbf{y})$ is therefore equal to $P(\mathbf{y})R_n(\mathbf{y})$. Now, assume that at step n the \mathbf{x} 's are actually distributed with the target distribution: then $p_n(\mathbf{x}) = P(\mathbf{x})$. It immediately follows that $R_n(\mathbf{y}) = 1$ and $p_{n+1}(\mathbf{y}) = P(\mathbf{y})$: we have thus proven that the target density is a fixed point of this algorithm.

9.4.5 An elementary case study for Gibbs

A very simple analysable case for Gibbs sampling is the two-dimensional density

$$P(\mathbf{x}) = x^1 + x^2 . \tag{194}$$

If at step n the density of walkers is given by $p_n(\mathbf{x})$ then after one Gibbs step we shall have (cf Eq.(192))

$$p_{n+1}(\mathbf{y}) = P(\mathbf{y}) \int dx^1 dx^2 \frac{p_n(\mathbf{x}) (y^1 + x^2)}{(x^2 + 1/2)(y^1 + 1/2)} \tag{195}$$

Let us start with a uniform density of walkers: $p_1(\mathbf{x}) = \theta(0 < x^{1,2} < 1)$. After the first Gibbs step, Eq.(195) yields

$$p_1(\mathbf{x}) = P(\mathbf{x}) \frac{2 - \tau + 2\tau x^1}{1 + 2x^1} , \quad \tau = \log(3) . \tag{196}$$

By direct computation we find that, in general

$$p_n(\mathbf{x}) = P(\mathbf{x}) \frac{a_n + b_n x^1}{1 + 2x^1} , \tag{197}$$

with a double recursion relation:

$$\begin{pmatrix} a_n \\ b_n \end{pmatrix} = \mathcal{M} \begin{pmatrix} a_{n-1} \\ b_{n-1} \end{pmatrix} , \quad \begin{pmatrix} a_1 \\ b_1 \end{pmatrix} = \begin{pmatrix} 2 - \tau \\ 2\tau \end{pmatrix} , \quad (198)$$

with

$$\mathcal{M} = \begin{pmatrix} 1 - \kappa & \kappa/2 \\ 2\kappa & 1 - \kappa \end{pmatrix} , \quad \kappa = \tau - \tau^2/2 = 0.49514 \dots \quad (199)$$

The matrix \mathcal{M} has eigenvalues 1 and $1 - 2\kappa = 0.0097 \dots$, for eigenvectors $v_1 = (1, 2)$ and $v_2 = (1, -2)$, respectively. Starting with the vector $(2 - \tau, 2\tau) = v_1 + (1 - \tau)v_2$ and repeatedly applying \mathcal{M} , we shall suppress the v_2 component exponentially fast, and find that $\lim_{n \rightarrow \infty} (a_n, b_n) = (1, 2)$, in other words,

$$\lim_{n \rightarrow \infty} p_n(\mathbf{x}) = P(\mathbf{x}) . \quad (200)$$

Since $1 - 2\kappa$ is so small, the approach to the target density is very fast in this case.

10 Phase space algorithms for particle physics

In particle physics, a recurrent integration problem is that of a differential cross section over the collection of allowed final states. The determination of the differential cross section is outside the scope of these lectures, rather we are interested in the phase space integration itself.

10.1 The uniform phase space problem in particle phenomenology

In particle phenomenology (*not* in statistical physics), *phase space* denotes the collection of all possible final-state *momenta*., including the constraints posed by the on-shell conditions on the individual momenta as well as the overall restriction posed by conservation of energy and momentum. Denoting the masses of the particles by m_j and the four-momenta by p_j^μ ($j = 1, 2, \dots, n$), and the *total* four-momentum by P^μ , the (relativistic) phase space integration element reads⁷⁶

$$dV_n(P; p_1, \dots, p_n) = \prod_{j=1}^n \left(d^4 p_j \delta(p_j^2 - m_j^2) \right) \delta^4 \left(P - \sum_{j=1}^n p_j \right) . \quad (201)$$

It defines a $(3n - 4)$ -dimensional subspace of the $4n$ -dimensional space of all momentum components. Because of the mass-shell conditions, this subspace is highly nonlinear. The ‘phase space problem’ is that of designing algorithms to generate phase space points (‘events’) with uniform probability. That this is far from trivial can be seen from the fact that the *volume* of phase space is in general unknown. The only⁷⁷ closed-form results are those for $n = 2$ with arbitrary masses, and for general n the *ultrarelativistic limit* with $m_j = 0$, and the *nonrelativistic limit* where $P^2 \approx (m_1 + \dots + m_n)^2$.

10.2 Two-body phase space

10.2.1 The two-body algorithm

The simplest phase space is the two-body one: by standard methods we have

$$dV_2(P; p_1, p_2) = d^4 p_1 \delta(p_1^2 - m_1^2) d^4 p_2 \delta(p_2^2 - m_2^2) \delta^4(P - p_1 - p_2)$$

⁷⁶There is an additional constraint that the energies of all momenta be *positive*. This is usually left to be understood.

⁷⁷As far as I know.

$$= \frac{1}{8} \mathcal{F} \left(\frac{m_1^2}{s}, \frac{m_2^2}{s} \right) d\Omega \ , \quad (202)$$

where $s = P^2$, Ω is the solid angle of \mathbf{p}_1 in the rest frame of P^μ , and

$$\mathcal{F}(x, y) = \left((1 - x - y)^2 - 4xy \right)^{1/2} . \quad (203)$$

The algorithm is given below; the momenta $p_{1,2}$ are generated in the rest frame of P , and then boosted to the frame in which P^μ was given. The

Algorithm 8 Two-body phase space with masses $m_{1,2}$ and total invariant energy $\sqrt{s} > m_1 + m_2$

$p_1^0 \leftarrow (s + m_1^2 - m_2^2)/2\sqrt{s}$
 $p_2^0 \leftarrow \sqrt{s} - p_1^0 \ , \quad q \leftarrow \sqrt{(p_1^0)^2 - m_1^2} \ \{\text{compute energies and momenta}\}$
 $\cos \theta \leftarrow -1 + 2\rho \ , \quad \phi \leftarrow 2\pi\rho \ \{\text{the solid angle}\}$
 $p_1^1 \leftarrow q \sin \theta \cos \phi \ , \quad p_1^2 \leftarrow q \sin \theta \sin \phi \ , \quad p_1^3 \leftarrow q \cos \theta$
 $\mathbf{p}_2 \leftarrow -\mathbf{p}_1 \ \{\text{construct the momenta in the } P \text{ rest frame}\}$
 boost $p_{1,2}^\mu$ to the actual frame of P^μ (see algorithm 9)

Lorentz boost is that which takes the vector $(\sqrt{s}, \vec{0})$ over into (P^0, \mathbf{P}) . We give it here separately.

Algorithm 9 Lorentz boost from P^μ , with $P^2 = s$, at rest to given form, applied on vector p^μ . The resultant vector is q^μ .

$q^0 \leftarrow (p^0 P^0 + p^1 P^1 + p^2 P^2 + p^3 P^3)/\sqrt{s} \ \{\text{Note signs!}\}$
 $\vec{q} \leftarrow \vec{p} + \vec{P} (p^0 + q^0)/(\sqrt{s} + P^0)$

10.2.2 Two-body reduction

Let us manipulate Eq.(201) by multiplying it with a clever choice of unity:

$$\begin{aligned}
 dV(P; p_1, p_2, \dots, p_n) &= \prod_{j=1}^n \left(d^4 p_j \delta(p_j^2 - m_j^2) \right) \delta^4 \left(P - \sum_{j=1}^n \right) \\
 &\quad \times d^4 q_1 \delta^4 \left(q_1 - \sum_{j=2}^n \right) du_1 \delta(q_1^2 - u_1) \\
 &= dV_2(P; p_1, q_1) du_1 dV_{n-1}(q_1; p_2, \dots, p_n) . \quad (204)
 \end{aligned}$$

We can continue this so that eventually

$$dV_n(P; p_1, \dots, p_n) = dV_2(P; p_1, q_1) du_1 dV_2(q_1; p_2, q_2) du_2 dV_2(q_2; p_3, q_3) du_3 \cdots \cdots dV_2(q_{n-3}; p_{n-2}, q_{n-2}) du_{n-3} dV_2(q_{n-2}; p_{n-1}, p_n) du_{n-2} \quad (205)$$

The n -body problem has now been split up into a cascade of $n - 1$ two-body problems and a selection of $n - 2$ invariant masses squared. Using the result for the two-body case, the *volume* of the n -body phase space is then given by an $(n - 2)$ -dimensional integral:

$$V_n(P) = \left(\frac{\pi}{2}\right)^{n-1} \int du_1 \cdots du_{n-2} \left(\prod_{j=1}^{n-2} \mathcal{F}\left(\frac{u_j}{u_{j-1}}, \frac{m_j^2}{u_{j-1}}\right) \right) \mathcal{F}\left(\frac{m_{n-1}^2}{u_{n-2}}, \frac{m_n^2}{u_{n-2}}\right) \quad (206)$$

where $u_0 = s$. The integration boundaries are

$$\sqrt{u_j} < \sqrt{u_{j-1}} - m_j \quad (j = 1, \dots, n-2) \quad , \quad \sqrt{u_{n-2}} > m_{n-1} + m_n \quad (207)$$

Unsurprisingly the general result for $V_n(P)$ is not known, nor do we have an algorithm for generating the u 's.

10.3 The relativistic problem

The phase space problem becomes simpler if we can neglect the masses, or at least assume that they are small(ish).

10.3.1 Two-body reduction algorithm

If we let all the particle masses vanish, Eq.(206) becomes simpler:

$$V_n^{(0)}(P) = \left(\frac{\pi}{2}\right)^{n-1} \int du_1 \cdots du_{n-2} \prod_{j=1}^{n-2} \left(1 - \frac{u_j}{u_{j-1}}\right) \quad (208)$$

with $s > u_1 > u_2 > \cdots > u_{n-2} > 0$. Following [31] we introduce new variables $v_j = u_j/u_{j-1}$ ($j = 1, \dots, n-2$) and then we find

$$\begin{aligned} V_n^{(0)}(P) &= \left(\frac{\pi}{2}\right)^{n-1} s^{n-2} \int_0^1 dv_1 \cdots dv_{n-2} \prod_{j=1}^{n-2} v_j^{n-2-j} (1 - v_j) \\ &= \left(\frac{\pi}{2}\right)^{n-1} \frac{s^{n-2}}{(n-1)!(n-2)!} \quad (209) \end{aligned}$$

The remaining problem is to generate variables v according to $v^k(1-v)$. It can be done by inversion, where we solve

$$\int_0^v dw w^k(1-w) = \rho \int_0^1 dw w^k(1-w) \quad \rightarrow \quad (k+2)v^{k+1} - (k+1)v^{k+2} = \rho \quad . \quad (210)$$

All this then leads to the following algorithm for massless n -body phase space.

Algorithm 10 The Platzzer algorithm for $n \geq 3$

generate v_j according to $(k+1)(k+2)v_j^{n-2-k}(1-v_j)$ ($j = 1, \dots, n-2$)
compute $u_j = v_j u_{j-1}$ ($j = 1, \dots, n-2$)
use algorithm 8 for the cascade $P \rightarrow p_1 + q_1$, $q_1 \rightarrow p_2 + q_2$, $q_2 \rightarrow p_3 + q_3$, \dots , $q_{n-2} \rightarrow p_{n-1} + p_n$

⌘ This way of handling many-body phase space has a long history; I have simply referred to its most recent incarnation.

10.3.2 Massless RAMBO

The algorithm of the previous section imposes a hierarchy on the momenta that has no physical basis, and involves $n-1$ Lorentz boosts than may lead to numerical inaccuracies [31]. On the other hand, we can generalise the Euler algorithm to impose the overall constraint of momentum conservation, and arrive at a more ‘democratic’ approach : the **RAMBO** algorithm that is another case of the Rule of Nifty. We start by generating unconstrained massless momenta, and then modify them to enforce the correct overall momentum. From

$$\int d^4q \delta(q^2) \exp(-q^0) = \int_0^\infty dq^0 \frac{q^0}{2} \exp(-q^0) \int d\Omega = 2\pi \quad (211)$$

we see that a UA reads

$$1 = \frac{1}{(2\pi)^n} \int \prod_{j=1}^n d^4q_j \delta(q_j^2) \exp(-q_j^0) \quad . \quad (212)$$

The momenta q_j add up to a total momentum Q . We then perform a scaling by a factor $\sqrt{Q^2/s}$ so that Q gets the right invariant mass, and a Lorentz boost Λ that takes the vector Q into its rest frame. We can write this out as

$$\begin{aligned}
1 &= \frac{1}{(2\pi)^n} \int \prod_{j=1}^n d^4 q_j \delta(q_j^2) \exp(-q_j^0) \\
&\quad d^4 Q \delta^4 \left(Q - \sum_{j=1}^n q_j \right) d(x^2) \delta \left(x^2 - \frac{Q^2}{s} \right) \\
&\quad \prod_{j=1}^n d^4 p_j \delta^4 \left(p_j - \frac{1}{x} \Lambda q_j \right) . \tag{213}
\end{aligned}$$

We now eliminate the q_j , using

$$\begin{aligned}
\delta^4 \left(\sum_{j=1}^n q_j - Q \right) &= \frac{1}{x^4} \delta \left(\sum_{j=1}^n p_j - \frac{1}{x} \Lambda Q \right) = \frac{1}{x^4} \delta^4 \left(\sum_{j=1}^n p_j - P \right) , \\
\delta^4 \left(p_j - \frac{1}{x} \Lambda q_j \right) &= x^4 \delta^4 (q_j - x \Lambda^{-1} p_j) , \quad \delta(q_j^2) = \frac{1}{x^2} \delta(p_j^2) , \tag{214}
\end{aligned}$$

so as to arrive at

$$\begin{aligned}
1 &= \mathcal{A} \int dV_n(P; p_1, \dots, p_n) , \\
\mathcal{A} &= \frac{1}{(2\pi)^n} \int d^4 Q d(x^2) \exp(-Q^0) x^{2n-4} \delta(x^2 - Q^2/s) \\
&= \frac{2(2\pi)^{n-1}}{s^{n-2}} \int_0^\infty dQ^0 \int_0^{Q^0} d|Q| |Q|^2 \left((Q^0)^2 - |Q|^2 \right)^{n-2} \exp(-Q^0) \\
&= 1/V_n^{(0)}(P) . \tag{215}
\end{aligned}$$

This proves the correctness of the **RAMBO** algorithm. Note that, as before, the damping factor $\exp(-q_j^0)$ precisely guarantees a uniform density. In addition we see that particle masses cannot be easily included because of the scaling by the *a priori* unknown scaling factor x . Finally, we see that $4n$ pseudo-random numbers are used to sample points in $(3n - 4)$ -dimensional space. Therefore, information is lost, and the $4n$ -dimensional space of random numbers contains subspaces of dimension $n + 4$ in which every point ends up in the same phase space point.

Algorithm 11 The RAMBO algorithm for n momenta with total invariant mass squared s

generate energies q_j^0 according to $q_j^0 \exp(-q_j^0)$, $1 \leq j \leq n$
generate momenta \mathbf{q}_j isotropically, with $|\mathbf{q}_j| = q_j^0$, $1 \leq j \leq n$
compute $Q^\mu = \sum_j q_j^\mu$, and $x = \sqrt{Q^2/s}$
 $p_j \leftarrow \Lambda q_j/x$, $1 \leq j \leq n$ $\{\Lambda$ is the Lorentz boost that brings Q to its rest frame, see algorithm12}
each event carries a weight given by Eq.(209)

Algorithm 12 Lorentz boost from P^μ , with $P^2 = s$, to rest from given form, applied on vector p^μ . The resultant vector is q^μ .

$q^0 \leftarrow (p^0 P^0 - p^1 P^1 - p^2 P^2 - p^3 P^3)/\sqrt{s}$ {Note signs!}
 $\vec{q} \leftarrow \vec{p} - \vec{P} (p^0 + q^0)/(\sqrt{s} + P^0)$

10.3.3 Inclusion of masses

For the case where the particle masses are nonzero no algorithm is known that populates phase space uniformly⁷⁸. However, if the masses (and n) are not too large the following procedure can be used. We start with generating massless momenta p_j . We then proceed to scale down the 3-momenta of the particles with a common factor ξ , which opens up room to increase the energies such that masses can be accommodated. We therefore set

$$\Phi(\xi) = \sum_{j=1}^n \sqrt{m_j^2 + \xi^2 |\mathbf{p}_j|^2} = \sqrt{s} \quad . \quad (216)$$

Since $\Phi(\xi)$ is monotonic and a solution with $0 < \xi < 1$ exists, it is not difficult to find the right value of ξ , and the UA description becomes

$$\begin{aligned} 1 &= \frac{1}{V_n(P)} \int \left(\prod_{j=1}^n d^4 p_j \delta(p_j^2 - m_j^2) \right) \delta^4 \left(\sum_{j=1}^n p_j - P \right) \\ &= \frac{1}{V_n(P)} \int \left(\prod_{j=1}^n \frac{1}{2|\mathbf{p}_j|} d^3 \mathbf{p}_j \right) \delta^3 \left(\sum_{j=1}^n \mathbf{p}_j \right) \delta \left(\sum_{j=1}^n |\mathbf{p}_j| - \sqrt{s} \right) \\ &\quad d\xi \Phi'(\xi) \delta(\Phi(\xi) - \sqrt{s}) \end{aligned}$$

⁷⁸Except for $n = 2$, see above, and for $n = 3$ with $m_1 > 0, m_{2,3} = 0$.

$$\left(\prod_{j=1}^n d^3 \mathbf{k}_j \delta^3(\mathbf{k}_j - \xi \mathbf{p}_j) dk_j^0 \delta(k_j^0 - \sqrt{|\mathbf{k}_j|^2 + m_j^2}) \right) , \quad (217)$$

which after standard⁷⁹ manipulations results in

$$\begin{aligned} 1 &= \int \left(\prod_{j=1}^n d^4 k_j \delta(k_j^2 - m_j^2) \right) \delta^4 \left(\sum_{j=1}^n k_j - P \right) \frac{\mathcal{G}}{V_n(P)} , \\ \mathcal{G} &= \left(\prod_{j=1}^n \frac{k_j^0}{|\mathbf{k}_j|} \right) \left(\sum_{j=1}^n \frac{|\mathbf{k}_j|^2}{k_j^0 \sqrt{s}} \right) \left(\sum_{j=1}^n \frac{|\mathbf{k}_j|}{\sqrt{s}} \right)^{3-2n} . \end{aligned} \quad (218)$$

The scaling algorithm therefore results in a nonuniform sampling, with event weights given by $V_n(P)/\mathcal{G}$. The scaling operation is reversible, and no information is lost⁸⁰. On the other hand, it is not known what the maximum weight is for general masses.

Algorithm 13 Giving masses to massless momenta

generate n massless momenta p_j using algorithm 11

find ξ such that $\sum_j \sqrt{|\mathbf{q}_j|^2 \xi^2 + m_j^2} = \sqrt{s}$

$\mathbf{k}_j \leftarrow \xi \mathbf{q}_j$, $k_j^0 \leftarrow \sqrt{|\mathbf{k}_j|^2 + m_j^2}$, $1 \leq j \leq n$

each event carries a weight given by $V_n(P)/\mathcal{G}$, see Eq.(218)

⁷⁹By now, hopefully!

⁸⁰We can reconstruct the ‘original’ p_j from the k_j .

⌘ Consider $n = 3$. In the massless case, the phase space (the *Dalitz plot*) is a triangle in terms of two of the energies: $0 < p_{1,2}^0 < \sqrt{s}/2$ and $p_1^0 + p_2^0 < \sqrt{s}$. When masses are introduced, the triangle is ‘contracted’ into a figure with rounded edges and no sharp points. A *uniform* mapping is therefore not possible, and the events will have nonconstant weights. Another message is that the ‘contraction’ may be expected to be minimal (and hence the weight maximal) for those events that are ‘as far as possible’ from the edges of phase space. This translates into the rule of thumb that the maximum weight is attained when all massless particles have zero energy, and all massive particles have the highest possible velocity (in whatever sense this makes sense).

10.4 Nonrelativistic phase space: BOLTZ

At the other extreme of the energy range we have nonrelativistic phase space. This also allows for a ‘democratic’ approach, as we shall now show. From

$$\int d^3\mathbf{q} \exp\left(-\frac{|\mathbf{q}|^2}{2m}\right) = (2\pi m)^{3/2} \quad (219)$$

we see that a UA can be given:

$$1 = A \int \prod_{j=1}^n d^3\mathbf{q}_j \exp\left(-\frac{|\mathbf{q}_j|^2}{2m}\right) \quad , \quad A = \prod_{j=1}^n (2\pi m_j)^{-3/2} \quad . \quad (220)$$

Obviously, no mass is allowed to vanish here. Having generated the \mathbf{q}_j we then perform a *Galilei* transformation that takes the momenta over into new momenta \mathbf{k}_j that add up to zero. Subsequently we scale the momenta \mathbf{k}_j into \mathbf{p}_j that have the correct total kinetic energy U . With the notation $\mathbf{Q} = \sum_j \mathbf{q}_j$, $M = \sum_j m_j$, $E_k = \sum_j |\mathbf{k}_j|^2/(2m_j)$, and $\hat{A} = (2\pi M)^{3/2} A$, we can analyse this procedure as follows, by successive elimination of the \mathbf{q}_j and the \mathbf{k}_j :

$$1 = A \int \left(\prod_{j=1}^n d^3\mathbf{q}_j \exp\left(-\frac{|\mathbf{q}_j|^2}{2m}\right) \right)$$

$$\begin{aligned}
& \times d^3\mathbf{v} \delta^3\left(\frac{\mathbf{Q}}{M} - \mathbf{v}\right) \left(\prod_{j=1}^n d^3\mathbf{k}_j \delta^3(\mathbf{k}_j - \mathbf{q}_j + m_j\mathbf{v})\right) \\
& = A \int \left(\prod_{j=1}^n d^3\mathbf{k}_j\right) \delta^3\left(\sum_{j=1}^n \mathbf{k}_j\right) d^3\mathbf{v} M^3 \exp\left(-E_k - \frac{M|\mathbf{v}|^2}{2}\right) \\
& = \hat{A} \int \left(\prod_{j=1}^n d^3\mathbf{k}_j\right) \delta^3\left(\sum_{j=1}^n \mathbf{k}_j\right) \exp(-E_k) \\
& \quad \times d(x^2) \delta\left(x^2 - \frac{E_k}{U}\right) \left(\prod_{j=1}^n d^3\mathbf{p}_j \delta^3\left(\mathbf{p}_j - \frac{1}{x}\mathbf{k}_j\right)\right) \\
& = \hat{A} \int \left(\prod_{j=1}^n d^3\mathbf{p}_j\right) \delta^3\left(\sum_{j=1}^n \mathbf{p}_j\right) \delta\left(\sum_{j=1}^n \frac{|\mathbf{p}_j|^2}{2m_j} - U\right) d(x^2) x^{3n-5} \exp(-Ux^2) U \\
& = \hat{A} \Gamma\left(\frac{3n-3}{2}\right) U^{(5-3n)/2} \int \left(\prod_{j=1}^n d^3\mathbf{p}_j\right) \delta^3\left(\sum_{j=1}^n \mathbf{p}_j\right) \delta\left(\sum_{j=1}^n \frac{|\mathbf{p}_j|^2}{2m_j} - U\right) .
\end{aligned} \tag{221}$$

We have thus proven that we can sample nonrelativistic phase space uniformly, and have also computed the volume of this phase space: it reads

$$V_{nr}(U; m_1, \dots, m_n) = \frac{\prod_{j=1}^n (2\pi m_j)^{3/2}}{(2\pi M)^{3/2}} \frac{U^{(3n-5)/2}}{\Gamma((3n-3)/2)} . \tag{222}$$

This Rule of Nifty result is the basis of the BOLTZ algorithm⁸¹.

Algorithm 14 The BOLTZ algorithm for total energy U and masses $m_{1,2,\dots,n}$

generate q_j^r according to $\exp(-(q_j^r)^2/(2m_j))$ for $r = x, y, z$ and $1 \leq j \leq n$
 $\mathbf{v} \leftarrow \sum_j \mathbf{q}_j / \sum_j m_j$
 $\mathbf{k}_j \leftarrow \mathbf{q}_j - m_j \mathbf{v}$ for $1 \leq j \leq n$
compute E_k and $x = \sqrt{E_k/U}$
 $\mathbf{p}_j \leftarrow \mathbf{k}_j/x$ for $1 \leq j \leq n$
each event carries a weight given by Eq.(222)

⌘ In trms of staistical physics, the —tt BOLTZ algorithm generates the *microcanonical* ensemble.

⁸¹Boltz, man!

⌘ I do not know of a spacetime transformation that ‘interpolates’ between a Galilei and a Lorentz transform. This makes it understandable why the general phase space problem is so hard.

11 Appendices

11.0.1 Falling powers

The falling powers are defined as

$$N^{\underline{k}} = N(N-1)(N-2)\cdots(N-k+1) = \frac{N!}{(N-k)!} \quad (223)$$

By its definition, we have $N^{\underline{k}} = 0$ when $k > N$, and $N^{\underline{N}} = N!$. For large N and finite k we can approximate

$$N^{\underline{k}} \approx N^k \left(1 - \frac{k(k-1)}{2N}\right) . \quad (224)$$

From $N^{\underline{k}} = N!/(N-k)!$ we find immediately the binomial sum

$$\sum_{k \geq 0} \frac{x^k}{k!} N^{\underline{k}} = \sum_{k \geq 0} \binom{N}{k} x^k = (1+x)^N . \quad (225)$$

We can extend this by summing over N as well:

$$\begin{aligned} \sum_{N \geq 0} \sum_{k \geq 0} \frac{x^k}{k!} N^{\underline{k}} y^N &= \sum_{N \geq 0} y^N (1+x)^N \\ &= \frac{1}{1-y-xy} = \sum_{m \geq 0} \frac{x^m y^m}{(1-y)^{m+1}} . \end{aligned} \quad (226)$$

Isolating from this the power x^k gives

$$\sum_{N \geq 0} y^N N^{\underline{k}} = \frac{k! y^k}{(1-y)^{k+1}} . \quad (227)$$

Eq.(225) also leads to the ‘binomial theorem for falling powers’:

$$(x+y)^{\underline{s}} = \sum_{r=0}^s \binom{s}{r} x^{\underline{r}} y^{\underline{s-r}} . \quad (228)$$

11.0.2 Relations between direct sums and unequal-sums

$$\begin{aligned}
S_1^2 &= S_2 + S_{1,1} , \\
S_2 S_1 &= S_3 + S_{2,1} , \\
S_1^3 &= S_3 + 3 S_{2,1} + S_{1,1,1} , \\
S_3 S_1 &= S_4 + S_{3,1} , \\
S_2^2 &= S_4 + S_{2,2} , \\
S_2 S_1^2 &= S_4 + 2 S_{3,1} + S_{2,2} + S_{2,1,1} , \\
S_1^4 &= S_4 + 4 S_{3,1} + 3 S_{2,2} + 6 S_{2,1,1} + S_{1,1,1,1} ;
\end{aligned} \tag{229}$$

and conversely,

$$\begin{aligned}
S_{1,1} &= S_1^2 - S_2 , \\
S_{2,1} &= S_2 S_1 - S_3 , \\
S_{1,1,1} &= S_1^3 - 3 S_2 S_1 + 2 S_3 , \\
S_{3,1} &= S_3 S_1 - S_4 , \\
S_{2,2} &= S_2^2 - S_4 , \\
S_{2,1,1} &= S_2 S_1^2 - 2 S_3 S_1 - S_2^2 + 2 S_4 , \\
S_{1,1,1,1} &= S_1^4 - 6 S_2 S_1^2 + 8 S_3 S_1 + 3 S_2^2 - 6 S_4 .
\end{aligned} \tag{230}$$

11.0.3 An exponential sum and the Poisson formula

Consider the function

$$\delta(s; x) = \sum_n s^{|n|} e^{2i\pi n x} = \frac{1 - s^2}{1 - 2s \cos(2\pi x) + s^2} , \tag{231}$$

where $0 < s < 1$. As $s \rightarrow 1$, $\delta(s; x)$ goes to zero except for integer x , where it approaches infinity. Moreover,

$$\int_{k-1/2}^{k+1/2} dx \delta(s; x) = 1 \tag{232}$$

for any integer k . Consequently we may write

$$\sum_n e^{2i\pi n x} = \sum_k \delta(x - k) \tag{233}$$

in the sense of distributions, that is integrated with a suitable test function. This implies the Poisson formula:

$$\sum_k f(k) = \sum_n g(n) \quad , \quad g(y) = \int dx f(x) e^{2i\pi xy} . \quad (234)$$

11.0.4 About the integral (40)

For M large, we can approximate

$$\begin{aligned} \log \left((1 + x/M)^M \right) &= M \log(1 + x/M) \\ &= M \left(\frac{x}{M} - \frac{x^2}{2M^2} + \frac{x^3}{3M^3} - \frac{x^4}{4M^4} + \cdots \right) \\ &= x - \frac{x^2}{2M} + \frac{x^3}{3M^2} - \frac{x^4}{4M^3} + \cdots \end{aligned} \quad (235)$$

We can therefore estimate the integral

$$\begin{aligned} \int_0^\infty dx e^{-x} \left(1 + \frac{x}{M} \right)^M &= \int_0^\infty dx \exp \left(-\frac{x^2}{2M} + \frac{x^3}{3M^2} - \frac{x^4}{4M^3} + \cdots \right) \\ &\approx \int_0^\infty dx e^{-x^2/2M} \left(1 + \frac{x^3}{3M^2} - \frac{x^4}{4M^3} + \frac{x^6}{18M^4} + \cdots \right) \\ &= \frac{1}{2} \sqrt{2\pi M} \left(1 - \frac{(3M^2)}{4M^3} + \frac{(15M^3)}{18M^4} \right) + \frac{(2M)^2}{6M^2} + \cdots \\ &= \sqrt{\frac{\pi M}{2}} \left(1 + \frac{1}{12M} \right) + \frac{2}{3} + \cdots \end{aligned} \quad (236)$$

The neglected terms are of order $1/M$ and smaller.

11.0.5 Selfies

Consider an ‘arbitrarily chosen algorithm’ of the type of Eq.(35). The number a is a *selfie* if $f(a) = a$. The probability that a given number is a selfie is $1/M$, and therefore the probability of having exactly k selfies in an arbitrarily chosen algorithm is

$$S_M(k) = \binom{M}{k} \left(\frac{1}{M} \right)^k \left(1 - \frac{1}{M} \right)^{M-k} . \quad (237)$$

Since

$$S_M(0) = \left(1 - \frac{1}{M}\right)^M \approx \frac{1}{e} \left(1 - \frac{1}{2M} + \mathcal{O}\left(\frac{1}{M^2}\right)\right) , \quad (238)$$

the probability to have *at least* one selfie is about $1 - 1/e \approx 63\%$. The largest probability is $S_M(1)$, and incidentally

$$\sum_{k=0}^M k S_M(k) = 1 , \quad (239)$$

so that the *expected number* of selfies is exactly 1, independent of M . This therefore also holds for shift-register PRNGs.

11.0.6 Serial correlation in a real-number model

Consider a multiplicative congruential PRNG with modulus m , multiplier a and increment c . The serial correlation is

$$r = \frac{\sum_{n=1}^k x_n x_{n+1} - \left(\sum_{n=1}^k x_n\right)^2}{\sum_{n=1}^k x_n^2 - \left(\sum_{n=1}^k x_n\right)^2} , \quad (240)$$

a version of the Pearson correlation coefficient. We can compute an approximate value for this correlation using the ‘real number’ model for the PRNG: after scaling by $1/m$, the $z_n = x_n/m$ values will be approximately the real numbers in $(0, 1)$. Then $z_{n+1} = y(z_n)$ with

$$y(z) = (az + \delta) \bmod 1 = (az_n + \delta) - \sum_{k=1}^a \theta\left(z > \frac{k - \delta}{a}\right) . \quad (241)$$

Then, assuming uniform distribution of the z_n values we will have approximately

$$\langle z_n^2 \rangle \approx 1/3 , \quad \langle z_n \rangle^2 \approx 1/4 , \quad (242)$$

and

$$\begin{aligned} \langle z_n y(z_n) \rangle &\approx \int_0^1 dz \, z y(z) = \int_0^1 dz \, (az^2 + \delta z) - \sum_{k=1}^a \int_{(k-\delta)/a}^1 dz \, z \\ &= -\frac{a}{6} + \frac{\delta}{2} - \frac{1}{2a^2} \sum_{k=1}^a (k^2 - 2\delta k + \delta^2) \\ &= \frac{1}{4} + \frac{1}{12a} - \frac{1}{2a} \delta(1 - \delta) . \end{aligned} \quad (243)$$

The estimate for the correlation is therefore

$$r \approx \frac{1}{a} \left(1 - 6\delta(1 - \delta) \right) . \quad (244)$$

11.0.7 The two-point function for the Euler diaphony

For the one-dimensional Euler diaphony we have

$$\beta(x) = \sum_{n \neq 0} \frac{3}{\pi^2 n^2} \exp(2i\pi n x) . \quad (245)$$

Now assume that $0 < x < 1$. By differentiating twice we have

$$\beta''(x) = -12 \sum_{n \neq 0} \exp(2i\pi n x) = 12 \quad (246)$$

by virtue of Eq.(233). Thus,

$$\beta(x) = 6x^2 + px + q \quad (247)$$

There are two conditions to be met:

$$\beta(x) = \beta(1 - x) \quad , \quad \int_0^1 dx \beta(x) = 0 . \quad (248)$$

This leads to

$$\beta(x) = 1 - 6x(1 - x) \quad , \quad 0 < x < 1 . \quad (249)$$

At every integer value of z the two-point function has a kink to make it periodic, so the final result must be

$$\beta(z) = 1 - 6\{x\} \left(1 - \{x\} \right) \quad , \quad \{x\} = x - \lfloor x \rfloor . \quad (250)$$

If we are interested only in $-1 \leq x \leq 1$ we may replace $\{x\}$ by $|z|$.

11.0.8 Rational denominators for continued fractions

The continued fraction with the smallest possible coefficients is

$$\phi_1 = [1, 1, 1, 1, 1, \dots] = \frac{1}{1 + \phi_1} = \frac{1}{2} (\sqrt{5} - 1) \quad , \quad (251)$$

the golden ratio; its approximant denominator obeys

$$q_n = \theta(n = 0, 1) + \theta(n \geq 2)(q_{n-1} + q_{n-2}) \ , \quad (252)$$

so that $q_n = F_n$ are the Fibonacci numbers 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, ... Fibonacci number. We can find the asymptotic behaviour by introducing the generating function

$$\begin{aligned} f_1(x) &= \sum_{n \geq 0} x^n q_n = 1 + x + x(f_1(x) - 1) + x^2 f_1(x) \\ &= \frac{1}{1 - x - x^2} \ . \end{aligned} \quad (253)$$

Its nearest singularity⁸² is at $x = \phi_1$. We can therefore approximate⁸³

$$f_1(x) \sim \frac{1}{(\phi_1 - x)(2\phi_1 + 1)} = \frac{1}{\phi_1(2\phi_1 + 1)} \frac{1}{1 - x/\phi_1} \ , \quad (254)$$

which gives⁸⁴

$$q_n \sim \frac{1}{\phi_1(2\phi_1 + 1)} \left(\frac{1}{\phi_1} \right)^{n+1} \sim (0.724) (1.618)^n \ . \quad (255)$$

A continued fraction with nonminimal coefficients is

$$\phi_2 = [2, 2, 2, 2, 2, \dots] = \frac{1}{2 + \phi_2} = -1 + \sqrt{2} \ . \quad (256)$$

Its approximant denominators, successively 1, 2, 5, 12, 29, 70, 169, 408, ... obey the recursion

$$q_n = \theta(n = 0) + \theta(n \geq 1)(2q_{n-1} + q_{n-2}) \ , \quad (257)$$

with generating function

$$f_2(x) = \sum_{n \geq 0} x^n q_n = \frac{1}{1 - 2x - x^2} \ . \quad (258)$$

⁸²That singularity that lie closest to the origin in the complex plane.

⁸³A polynomial $P(x)$ with nondegenerate roots x_j is approximated by $(x - x_j)P'(x_j)$ close to the root.

⁸⁴The *exact* result also contains powers of $1/x_1$ that are therefore exponentially suppressed: the approximation is therefore very good already for moderate n .

Its nearest singularity is at $x = \phi_2$ and proceeding as before we find

$$q_n = \frac{1}{\phi_2(2\phi_2 + 2)} \left(\frac{1}{\phi_2} \right)^n \sim (0.854) (2.414)^n . \quad (259)$$

A continued fraction with ‘quite small’ coefficients is for instance

$$\phi_{12} = [1, 2, 1, 2, 1, 2, \dots] = \frac{1}{1 + \frac{1}{2 + \phi_{12}}} = \frac{2 + \phi_{12}}{3 + \phi_{12}} = -1 + \sqrt{3} , \quad (260)$$

and for its approximant denominator we have

$$\begin{aligned} &= \theta(n = 0, 1) + \theta(n \geq 2, \text{even})(2q_{n-1} + q_{n-2}) + \theta(n \geq 2, \text{odd})(q_{n-1} + q_{n-2}) \\ &= \theta(n = 0, 1) + \theta(n \geq 2)(q_{n-1} + q_{n-2}) + \theta(n \geq 2, \text{even})q_{n-1} , \end{aligned} \quad (261)$$

and for the generating function we find

$$f_{12}(x) = \sum_{n \geq 0} x^n q_n = 1 + x f_{12}(x) + x^2 f_{12}(x) + x \sum_{n \geq 1} x^{2n-1} q_{2n-1} . \quad (262)$$

Splitting $f_{12}(x)$ into even and odd parts we find that there must be an even function $g(x)$ such that $f_{12}(x) = (1 + x - x^2)g(x)$, and we finally arrive at

$$f_{12}(x) = \frac{1 + x - x^2}{(1 - x^2)^2 - 2x^2} . \quad (263)$$

The resulting denominators are, successively, 1, 1, 3, 4, 11, 15, 41, 56, 153, ..., growing faster than the Fibonacci numbers. The generating function has its nearest singularities at $\pm y_0$, with $y_0 = (\sqrt{3} - 1)/\sqrt{2}$, and we can therefore approximate

$$\begin{aligned} q_n &\sim \frac{1}{4\sqrt{3}} [(\sqrt{2} + 1) + (-1)^n(\sqrt{2} - 1)] \left(\frac{1}{y_0} \right)^{n+1} \\ &= \frac{\sqrt{3} + 1}{\sqrt{24}} \left(\frac{1}{y_0} \right)^n [\theta(n \text{ odd}) + \sqrt{2} \theta(n \text{ even})] \\ &\sim (0.789) (1.932)^n , n \text{ even} , (0.558) (1.932)^n , n \text{ odd} . \end{aligned} \quad (264)$$

The growth rate of the denominators is indeed inbetween that for ϕ_1 and ϕ_2 .

References

- [1] W.H. Press, B.P. Flannery, S.A. Teukolsky and W.T. Vetterling, *Numerical Recipes: The Art of Scientific Computing*, Cambridge University Press (1986).
- [2] *Nineteenth Century Clouds over the Dynamical Theory of Heat and Light*, by the Right. Hon. Lord Kelvin, G.C.V.O., D.C.L., LL.D., F.R.S, M.R.I, *Philosophical Magazine* S. 6. Vol. **2**, no. 7. July 1901.
- [3] N. Metropolis, S. Ulam, *The Monte Carlo Method*, J. Am. Stat. Ass. **44** (1949) 335.
- [4] *Monte Carlo Method* (Proceedings of a symposium in Los Angeles CA, sponsored by the RAND Corporation, the National Bureau of Standards, and the Oak Ridge National Laboratory, June 29 - July 1, 1949), National Bureau of Standards Applied Mathematics Series **12** (1951).
- [5] T.F. Chan, G.H. Golub, R.J. Leveque, *Algorithms for Computing the Sample Variance: Analysis and Recommendations*, The American Statistician **37** (1983) 242.
- [6] R. Bakx, R. Kleiss, F. Versteegen, *First- and second-order error estimates in Monte Carlo integration*, Comp. Phys. Comm. **208** (2016) 29.
- [7] N. Eldredge, S. J. Gould , *Punctuated equilibria: an alternative to phyletic gradualism*, in T.J.M. Schopf, ed., *Models in Paleobiology*. (San Francisco: Freeman, Cooper and Company), pp. 82?115.
- [8] D.E. Knuth, *The Art of Computer Programming, vol.2: Seminumerical Algorithms*, Addison-Wesley, 1981.
- [9] T.E. Hull, A.R. Dobell, *Random Number Generators* SIAM Review. **4**(3) 230.
- [10] R.D. Carmichael, *Note on a new number theory function* Bull. Am.Math.Soc. **16**(1910)232.
- [11] G. Marsaglia, *Random Numbers Fall Mainly in the Planes*, Proc. Nat. Acad. Sci. U.S.A. **61**(1)25.

- [12] System/360 Scientific Subroutine Package, Version III, Programmer's Manual. IBM, White Plains, New York, 1968, p. 77.
- [13] <http://random.mat.sbg.ac.at/results/karl/server/server.html>
- [14] M. Lüscher, *A Portable high quality random number generator for lattice field theory simulations*, Comput.Phys.Commun. **79**(1994) 100. This is an adaptation from the original algorithm proposed in G. Marsaglia and A. Zaman, *A New Class of Random Number Generators*, Ann. Appl. Prob. **1** (1991)462.
- [15] M. Matsumoto and T. Nishimura, *Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-random number generator*, ACM Trans. Mod. Comp. Sim. **8** (1998).
- [16] H. Leeb, MSc thesis, University of Salzburg, 1995.
- [17] L. Kuipers and H. Niederreiter, *Uniform Distribution of Sequences*, (Dover, 2006).
- [18] H. Woźniakowski, *Average case complexity of multivariate integration*, Bull. AMS **24** (1991) 185.
- [19] F. James, J. Hoogland, R. Kleiss, *Multidimensional sampling for simulation and integration: Measures, discrepancies, and quasirandom numbers*, Comp. Phys. Comm. **99** (1997) 180.
- [20] A. van Hameren, R. Kleiss, J. Hoogland, *Gaussian limits for discrepancies*, Nucl.Phys.Proc.Suppl. **63** (1998) 988-990.
- [21] J.M. Hammersley and D.C. Handscomb, *Monte Carlo Methods* (London: Methuen 1964).
- [22] B. Bratley, P. Fox and H. Niederreiter. AMC Transactions on Modeling and Computer Simulation 2, 3:195.
- [23] A. Lazopoulos, PhD thesis, Radboud University, Nijmegen (2007).
- [24] L. Devroye, *Non-uniform Random Variate Generation* (Springer, New York 1986). Out of print, it is now available for free on www.nrbook.com/devroye.

- [25] J. von Neumann, *Various Techniques Used in Connection With Random Digits*, contribution to [4]. This 3-page paper is an absolute must for anyone interested in Monte Carlo and its history.
- [26] M. Abramowitz and I. Stegun (eds.), *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, (Nat. Bur. Standards, 1964).
- [27] G.E.P. Box, M.E. Muller, *A Note on the Generation of Random Normal Deviates.*, Ann. Math. Stat. **29 (2)** (1958) 610.
- [28] A.J. Kinderman, J.F. Monahan, *Computer Generation of Random Variables Using the Ratio of Uniform Deviates*, ACM Trans. Math. Softw., **3 (3)** (1977).
- [29] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, E. Teller, *Equations of State Calculations by Fast Computing Machines*, J. Chem. Phys. **21 (6)** (1953) 1087;
W.K. Hastings, *Monte Carlo Sampling Methods Using Markov Chains and Their Applications*, Biometrika. **57 (1)** (1970): 97.
- [30] S.D. Ellis, R. Kleiss, W.J. Stirling, *A New Monte Carlo Treatment of Multiparticle Phase Space at High energies*. Comp. Phys. Comm. **40** (1986) 359
- [31] S. Platzer, *RAMBO on diet*, arXiv:1308.2922

Index

- billets, [10](#)
- Boltzmann's constant, [96](#)
- Buffon, [7](#)

- canonical ensemble, [96](#)
- continued-fraction representation, [73](#)
- contracted, [107](#)

- democracy, [103](#)
- detailed balance, [93](#)

- fixed point, [94](#)

- good lattice points, [71](#)
- gradient, [70](#)

- Ising system, [96](#)

- Lord Kelvin, [10](#)

- MCMC, [93](#)
- microconfigurations, [96](#)
- Monte Carlo, [11](#)
 - estimators, [18](#)
 - improved, [20](#)
 - numerical stability, [20](#)
 - positivity, [19](#)
 - integration, [11](#)
 - simulation, [11](#)
- most irrational number, [74](#)

- nearest singularity, [115](#)
- nonrelativistic limit, [100](#)

- phase space problem, [100](#)
- PRNG, [29](#)
 - logistic-map, [34](#)
 - Mersenne Twister, [39](#)
 - Midsquare, [32](#)
 - RANDU, [37](#)
 - RCARRY, [37](#)
 - shift-register, [32](#)
- programming consultant, [1](#)
- QRNG, [72](#)
- quadratic irrational, [74](#)

- Random numbers, [11](#)
 - probability definition, [11](#)
 - stream, [11](#)
- Rule of Nifty, [87](#)

- simulated annealing, [96](#)

- target density, [93](#)
- Temple column, [7](#)
- TGFSRPRNG, [39](#)
- Travelling Salesman, [96](#)

- ultrarelativistic limit, [100](#)