

# MA3227 Numerical Analysis II

## Lecture 22: Monte Carlo Methods

Simon Etter



2019/2020

# Monte Carlo Methods

## **Introduction**

Monte Carlo refers to a broad class of algorithms which use randomness to compute a deterministic outcome.

These methods and their name originated from the Manhattan project when the United States developed the first atomic bombs. Since this work was secret, the scientists working on the project required a code name for the methods that they were using. They settled on Monte Carlo in reference to a casino of the same name.

Monte Carlo methods tend to be quite slow, but they are a useful method of last resort when more traditional algorithms break down.

The example on the following slides will help make this clearer.

# Monte Carlo Methods

## Introductory example

Assume we want to compute the integral

$$I := \int_0^1 \dots \int_0^1 f(y_1, \dots, y_d) dy_1 \dots dy_d.$$

Applying a one-dimensional quadrature rule  $(x_k, w_k)_{k=1}^n$  with error  $e_n = \mathcal{O}(n^{-p})$  repeatedly, we obtain

$$\begin{aligned} I &= \int_0^1 \dots \int_0^1 \left( \sum_{k_1=1}^n f(x_{k_1}, y_2, \dots, y_d) w_{k_1} + \mathcal{O}(n^{-p}) \right) dy_2 \dots dy_d \\ &= \dots \\ &= \sum_{k_d=1}^n \dots \sum_{k_1=1}^n f(x_{k_1}, \dots, x_{k_d}) w_{k_1} \dots w_{k_d} + \mathcal{O}(n^{-p}). \end{aligned}$$

Observation: we need  $N = n^d$  function evaluations (all combinations of  $k_1, \dots, k_d \in \{1, \dots, n\}$ ) to achieve a  $\mathcal{O}(n^{-p}) = \mathcal{O}(N^{-p/d})$  error, i.e. the order of convergence in terms of  $N$  decreases for increasing  $d$ .

# Monte Carlo Methods

## Curse of dimensionality

The observation on the previous slide applies to a vast number of algorithms: the accuracy scales with the “effort per dimension”  $n$ , but the overall effort scales with  $n^d$ .

This phenomenon is called the curse of dimensionality and often makes it virtually impossible to do computations in dimensions  $d$  larger than about 4 or 5.

Problems with dimensions  $d > 5$  are more common than one might think. For example, the probability that someone has the corona virus is a function

$$\underbrace{\mathbb{R}}_{\text{body temperature}} \times \underbrace{\{0, 1\}}_{\text{cough}} \times \underbrace{\{0, 1\}}_{\text{runny nose}} \times \underbrace{\{0, 1\}}_{\text{travel history}} \times \underbrace{\{1, \dots, 5\}}_{\text{ethnicity}} \times \underbrace{\{1, \dots, 193\}}_{\text{nationality}} \rightarrow [0, 1].$$

Assuming we need 40 points to cover a reasonable range of body temperatures, this means we need to store

$$40 \times 2 \times 2 \times 2 \times 5 \times 193 = 308'800$$

numbers to represent this function. Moreover, each additional variable in our model multiplies the memory requirements.

# Monte Carlo Methods

## Introductory example (continued)

The curse of dimensionality can be overcome by reinterpreting  $I$  as the expectation value

$$I = \mathbb{E}[f(X_1, \dots, X_d)], \quad X_i \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1].$$

I will define the precise meaning of  $X_i \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$  later.

Expectations can be computed by taking the average of a sufficiently large number of samples  $x_i^{(k)}$  of the  $X_i$ , i.e.

$$\mathbb{E}[f(X_1, \dots, X_d)] = \frac{1}{N} \sum_{k=1}^N f(x_1^{(k)}, \dots, x_d^{(k)}) + \mathcal{O}(N^{-1/2}).$$

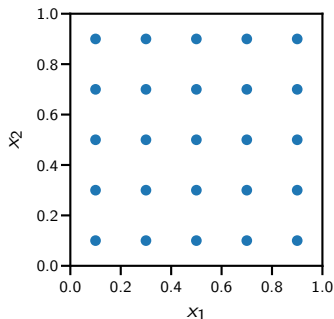
We will discuss the  $\mathcal{O}(N^{-1/2})$  error estimate later.

Observation:  $N$  evaluations of  $f$  lead to  $\mathcal{O}(N^{-1/2})$  error regardless of  $d$ . In particular, Monte Carlo is better than midpoint rule ( $p = 2$ ) for  $d > 4$ . See `convergence()` for numerical demonstration, and see the next slide for a pictorial illustration.

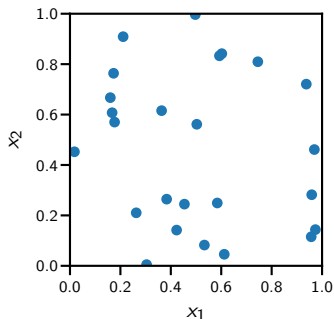
# Monte Carlo Methods

## Introductory example (continued)

Distribution of the function evaluation points  $(x_1, x_2)$ .



Midpoint



Monte Carlo

# Monte Carlo Methods

## Monte Carlo methods

In abstract terms, Monte Carlo corresponds to the approximation

$$\mathbb{E}[F] \approx \frac{1}{N} \sum_{k=1}^N F_k \quad \text{where} \quad F, F_k \stackrel{\text{iid}}{\sim} \mathcal{F}, \quad (1)$$

i.e. it amounts to estimating expectations by taking the mean of a sufficiently large number of samples as illustrated on the previous slide.

## Outlook

The aim of the following slides is to clarify and explain various technical details in (1) by revisiting some key definitions and results from probability theory.

I strongly recommend that you treat the following simply as an abstract mathematical theory in terms of sets and functions at first. Of course, our aim will be to eventually use this theory to understand why and how the Monte Carlo method works, but that discussion will be easier once we have a firm grasp on the underlying mathematics.

# Monte Carlo Methods

## **Disclaimer**

As usual in this module, I will take an engineer's approach to mathematics, i.e. I will ignore many technical details like measurability or  $\sigma$ -algebras since they are hardly ever relevant in applications.



# Monte Carlo Methods

## Def: Measure on a set $\Omega$

Function  $P$  mapping subsets of  $\Omega$  to nonnegative real numbers such that

$$P(\{\}) = 0 \quad \text{and} \quad P(A \cup B) = P(A) + P(B) \text{ if } A, B \subset \Omega \text{ are disjoint.}$$

## Def: Probability measure on a set $\Omega$

A measure  $P$  such that  $P(\Omega) = 1$ .

## Def: Probability space

A pair  $(\Omega, P)$  where  $P$  is a probability measure on  $\Omega$

## Def: Random variables

Any function  $X : \Omega \rightarrow \Xi$  defined on a probability space  $(\Omega, P)$ .

## Def: Distribution of a random variable $X : \Omega \rightarrow \Xi$

The probability measure  $\hat{P}$  on  $\Xi$  given by  $\hat{P}(A) = P(X^{-1}(A))$ .

It is common practice to write  $P(X \in A)$  instead of  $P(X^{-1}(A))$ , and  $P(X_1 \in A_1, X_2 \in A_2)$  instead of  $P(X_1^{-1}(A_1) \cap X_2^{-1}(A_2))$ .

This can be confusing if we forget that  $P(X \in A)$  is just a shorthand for  $P(X^{-1}(A))$ , so try your best not to forget this.

# Monte Carlo Methods

## Discussion

It is possible and common practise to talk about random variables and their distributions without ever specifying what the underlying probability space is.

For example, when we say  $X \sim \text{Uniform}[0, 1]$  (“ $X$  is uniformly distributed in  $[0, 1]$ ”), we mean that  $X : \Omega \rightarrow [0, 1]$  is a function defined on some unspecified probability space  $\Omega$  such that

$$\begin{aligned} P(X \in [a, b]) &= P(X^{-1}([a, b])) \\ &= P(\{\omega \in \Omega \mid X(\omega) \in [a, b]\}) \\ &= b - a \end{aligned}$$

for all  $a, b$  such that  $0 \leq a \leq b \leq 1$ .

There are countless ways how we could construct such an  $X$ . The simplest construction is to choose  $\Omega = [0, 1]$ ,  $P([a, b]) = b - a$  and  $X(\omega) = \omega$ , but equally well we could also set  $X(\omega) = 1 - \omega$ , or  $\Omega = [0, 2]$ ,  $P([a, b]) = \frac{b-a}{2}$ ,  $X(\omega) = \frac{\omega}{2}$ , or an even more complicated construction.

The point is that when we say  $X \sim \text{Uniform}[0, 1]$ , we do not care how  $X$  comes about as long as  $P(X \in [a, b]) = b - a$ .

# Monte Carlo Methods

## Representations of distributions

Recall that “distribution of a random variable  $X : \Omega \rightarrow \Xi$ ” refers to the probability measure  $\hat{P}(A) = P(X^{-1}(A)) = P(X \in A)$  on  $\Xi$ .

Furthermore, recall that a measure on  $\Xi$  is just a function mapping subsets of  $\Xi$  to nonnegative real numbers.

There are many ways to represent measures, and correspondingly there are many ways to represent distributions.

The three most common ways to represent distributions are listed on the next slide.

# Monte Carlo Methods

## Representations of distributions (continued)

- ▶ Probability mass function (PMF) of  $X : \Omega \rightarrow \Xi$  where  $\Xi$  is discrete: a function  $p : \Xi \rightarrow [0, 1]$  such that for all  $A \subset \Xi$

$$P(X \in A) = \sum_{x \in A} p(x).$$

- ▶ Probability density function (PDF) of  $X : \Omega \rightarrow \mathbb{R}^n$ : a function  $f : \mathbb{R}^n \rightarrow [0, \infty)$  such that for all  $A \subset \mathbb{R}^n$

$$P(X \in A) = \int_A f(x) dx.$$

- ▶ Cumulative distribution function (CDF) of  $X : \Omega \rightarrow \mathbb{R}$ : a function  $F : \mathbb{R} \rightarrow [0, 1]$  such that for all  $a \in \mathbb{R}$

$$P(X \leq a) = F(a).$$

# Monte Carlo Methods

## Example

	Uniform[0, 1]	$\mathcal{N}(\mu, \sigma^2)$
Name	Uniform distribution on [0, 1].	Normal distribution with mean $\mu$ and standard deviation $\sigma$
PDF	$f(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases}$	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$
CDF	$F(x) = \begin{cases} 0 & \text{if } x < 0, \\ x & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x > 1. \end{cases}$	$F(x) = \int_{-\infty}^x f(x) dx$ (No explicit formula available)

## Example

$$X \sim \text{Bernoulli}(p) \quad \Longleftrightarrow \quad P(X = 0) = 1 - p, \quad P(X = 1) = p.$$

# Monte Carlo Methods

## Def: Expectation

The expectation of a random variable  $X : \Omega \rightarrow \mathbb{R}^n$  is given by

$$\mathbb{E}[X] = \int X(\omega) d\omega$$

where  $\int d\omega$  denotes the integral with respect to the probability measure  $P$  of the underlying probability space.

Integration with respect to a measure may not have been introduced to you yet. If so, ignore the above and study the below special cases instead.

Again, it is possible and common practice to work with expectations without referring to the underlying probability space  $(\Omega, P)$ .

For example, one can show that if  $X : \Omega \rightarrow \Xi$  with  $\Xi$  discrete or  $\Xi = \mathbb{R}^n$  has probability mass / density  $p(x)$ , then

$$\mathbb{E}[X] = \sum_{x \in \Xi} x p(x) \quad \text{or} \quad \mathbb{E}[X] = \int_{\mathbb{R}^n} x p(x) dx,$$

and similarly for any function  $f$  on  $\Xi$  we have

$$\mathbb{E}[f(X)] = \sum_{x \in \Xi} f(x) p(x) \quad \text{or} \quad \mathbb{E}[X] = \int_{\mathbb{R}^n} f(x) p(x) dx.$$

# Monte Carlo Methods

## Def: Variance

The variance of a random variable  $X : \Omega \rightarrow \mathbb{R}^n$  is given by

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2].$$

## Useful formulae

Let  $X, Y : \Omega \rightarrow \mathbb{R}^n$  be random variables and  $a, b \in \mathbb{R}$ . Then,

- ▶  $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ .
- ▶  $\text{Var}[aX + b] = a^2 \text{Var}[X]$ .
- ▶  $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

# Monte Carlo Methods

## Def: Independence

A collection of random variables  $(X_i : \Omega \rightarrow \Xi_i)_{i=1}^n$  are called independent if for all  $(A_i \subset \Xi_i)_{i=1}^n$  it holds

$$P(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n P(X_i \in A_i).$$

One can show that  $(X_i)_{i=1}^n$  are independent if and only if

- ▶  $f(x_1, \dots, x_n) = f_1(x_1) \dots f_n(x_n)$  for all  $(x_i)_{i=1}^n$ , where  $f, f_i$  are the PDFs / PMFs of  $(X_i)_{i=1}^n$  and  $X_i$ , respectively.
- ▶  $F(x_1, \dots, x_n) = F_1(x_1) \dots F_n(x_n)$  for all  $(x_i)_{i=1}^n$ , where  $F, F_i$  are the CDFs of  $(X_i)_{i=1}^n$  and  $X_i$ , respectively.

Moreover, if  $(X_i)_{i=1}^n$  are independent, then

- ▶  $\mathbb{E}(X_1 \dots X_n) = \prod_{i=1}^n \mathbb{E}(X_i)$ .
- ▶  $(f_i(X_i))_{i=1}^n$  are independent for any collection of functions  $f_i$ .



# Monte Carlo Methods

## **Def: Independently and identically distributed**

We say a collection of random variables  $X_1, \dots, X_n$  is identically and independently distributed according to a distribution  $\mathcal{D}$ , or  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{D}$ , if the  $X_i$  are independent and each  $X_i \sim \mathcal{D}$ .

## **Def: Monte Carlo estimate of the expectation**

Assume  $F$  is a random variable with distribution  $\mathcal{F}$ . We then set

$$\tilde{\mathbb{E}}_N[F] = \frac{1}{N} \sum_{k=1}^N F_k \quad \text{where } F_k \stackrel{\text{iid}}{\sim} \mathcal{F}.$$

# Monte Carlo Methods

## Discussion

If you followed my advice on slide 7, then you think of the above simply as a collection of abstract definitions in terms of sets and functions. In particular, the Monte Carlo estimate  $\tilde{\mathbb{E}}_N[F]$  as defined above is just a function  $\Omega \rightarrow \Xi$  since it is the sum of functions  $F_k : \Omega \rightarrow \Xi$ .

Of course, the goal of a Monte Carlo simulation is not to produce a function  $\tilde{\mathbb{E}}_N[F]$  but rather to determine a single element of  $\Xi$  which is close to  $\mathbb{E}[F]$ . This can be achieved by choosing some fixed  $\omega \in \Omega$  and evaluating  $\tilde{\mathbb{E}}_N[F](\omega)$ .

The reason why this procedure works is because the central limit theorem presented on the next slide asserts that  $\tilde{\mathbb{E}}_N[F]$  maps almost all  $\omega \in \Omega$  to points  $\tilde{\mathbb{E}}_N[F](\omega)$  which are close to  $\mathbb{E}[F]$ .

# Monte Carlo Methods

## Central limit theorem

Assume  $F, F_1, \dots, F_N \stackrel{\text{iid}}{\sim} \mathcal{F}$ . Then,

$$\tilde{\mathbb{E}}_N[F] \xrightarrow{d} \mathcal{N}\left(\mathbb{E}[F], \frac{1}{N}\text{Var}[F]\right) \quad \text{for } N \rightarrow \infty.$$

Here,  $\xrightarrow{d}$  stands for convergence in distribution. This means that the distribution of the random variable on the left becomes increasingly indistinguishable from a normal distribution with mean and standard deviation as indicated on the right.

The above formulation of the central limit theorem is not quite correct because the right-hand side still involves  $N$  and is therefore not a valid limit. The technically correct formulation is

$$\sqrt{\frac{N}{\text{Var}[F]}} \left( \tilde{\mathbb{E}}_N[F] - \mathbb{E}[F] \right) \xrightarrow{d} \mathcal{N}(0, 1).$$

In this module, we will mostly ignore the convergence part and simply assume that

$$\frac{1}{N} \sum_{k=1}^N F_k \sim \mathcal{N}\left(\mathbb{E}[F], \frac{1}{N}\text{Var}[F]\right) \quad \text{if } N \gtrsim 100.$$

# Monte Carlo Methods

## Interpretation of the central limit theorem

Two slides ago, I claimed that the central limit theorem implies that  $\tilde{\mathbb{E}}_N[F]$  maps almost all  $\omega \in \Omega$  to points  $\tilde{\mathbb{E}}_N[F](\omega)$  close to  $\mathbb{E}[F]$ . There are at least two ways how we can see that this is true.

*Interpretation 1.* The central limit theorem in particular implies that for  $N \gtrsim 100$  we have

$$\mathbb{E}[\tilde{\mathbb{E}}_N[F]] = \mathbb{E}[F] \quad \text{and} \quad \text{Var}[\tilde{\mathbb{E}}_N[F]] = \frac{1}{N} \text{Var}[F].$$

See the remark at the bottom of the previous slide regarding the origin and implication of the assumption  $N \gtrsim 100$ .

Furthermore, recall that the variance is defined as

$$\text{Var}[\tilde{\mathbb{E}}_N[F]] = \mathbb{E}\left[\left(\tilde{\mathbb{E}}_N[F] - \mathbb{E}[\tilde{\mathbb{E}}_N[F]]\right)^2\right].$$

Inserting the former equations into the latter yields

$$\mathbb{E}\left[\left(\tilde{\mathbb{E}}_N[F] - \mathbb{E}[F]\right)^2\right] = \frac{1}{N} \text{Var}[F],$$

i.e.  $\tilde{\mathbb{E}}_N[F]$  has an error  $\sqrt{\frac{1}{N} \text{Var}[F]}$  in expectation.

# Monte Carlo Methods

## Interpretation of the central limit theorem (continued)

*Interpretation 2.* If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then

$$\begin{aligned} P(|X - \mu| > \varepsilon \sigma) &= \frac{1}{\sqrt{2\pi} \sigma} \int_{\mu - \varepsilon \sigma}^{\mu + \varepsilon \sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\varepsilon}^{\varepsilon} \exp\left(-\frac{y^2}{2}\right) dy \\ &= g(\varepsilon) \end{aligned}$$

where the second line follows by substituting  $y = \frac{x - \mu}{\sigma}$ ,  $dy = \frac{dx}{\sigma}$ .  
Moreover, it can be shown that  $g(\varepsilon) \rightarrow 1$  very quickly for  $\varepsilon \rightarrow \infty$ .

Since the central limit theorem asserts  $\tilde{\mathbb{E}}_N[F] \sim \mathcal{N}(\mathbb{E}[F], \frac{1}{N} \text{Var}[F])$ , we conclude that it is very unlikely that the Monte Carlo estimate  $\tilde{\mathbb{E}}_N[F]$  assumes values further than three or four standard deviations  $\sqrt{\frac{1}{N} \text{Var}[F]}$  from the exact value  $\mathbb{E}[F]$ .

# Monte Carlo Methods

## Discussion

Both of the above interpretations of the central limit theorem lead us to the conclusion that  $\tilde{\mathbb{E}}_N[F]$  acts like a funnel: it takes in most (but not necessarily all, see below) of  $\Omega$  and maps it into a region around  $\mathbb{E}[F]$  with a diameter of roughly  $\frac{1}{N}\text{Var}[F]$ .

However, the central limit theorem tell us nothing about the worst-case behaviour. It is possible that there are  $\omega$  such that  $\tilde{\mathbb{E}}_N[F](\omega)$  does not converge to  $\mathbb{E}[F]$  regardless of how many samples  $N$  that we take.

In a sense, Monte Carlo methods are therefore indeed like gambling: the success or failure of our computations is entirely dependent on whether or not we happen to choose a “good”  $\omega$ .

# Monte Carlo Methods

## Discussion (continued)

However, unlike real gambling, Monte Carlo methods allow us to stack the odds heavily in our favour. For example, if we take  $N = 10'000$  samples of a random variable  $F$  with  $\text{Var}[F] = \sigma^2$ , then

$$\begin{aligned} P(|\tilde{\mathbb{E}}_N[F] - \mathbb{E}[F]| > 0.1 \sigma) &= P\left(|\tilde{\mathbb{E}}_N[F] - \mathbb{E}[F]| > 0.1 \sqrt{N} \left(\frac{1}{\sqrt{N}} \sigma\right)\right) \\ &= 1 - g(0.1 \sqrt{N}) \approx 10^{-45}. \end{aligned}$$

See slide 21 regarding  $g(x)$ .

Such odds for picking a bad  $\omega$  are negligible compared to those of other sources of error (e.g. faulty hardware, programming errors, getting struck by lightning, etc.).

This concludes our discussion of the probability theory underlying Monte Carlo algorithms. The next lecture will show you how to translate this theory into actual code.

# Monte Carlo Methods

## Summary

- ▶ Monte Carlo idea: estimate expectations by taking the mean of a large number of samples.
- ▶ Basic probability theory: probability space, random variables, distributions, expectation, variance, independence.
- ▶ Central limit theorem: if  $F, F_1, \dots, F_N \stackrel{\text{iid}}{\sim} \mathcal{F}$ , then

$$\frac{1}{N} \sum_{k=1}^N F_k \xrightarrow{d} \mathcal{N}\left(\mathbb{E}[F], \frac{1}{N} \text{Var}[F]\right) \quad \text{for } N \rightarrow \infty.$$