

Statistical analysis of experimental data

Markov Chains

Aleksander Filip Żarnecki



Lecture 13

January 26, 2023

Markov Chains

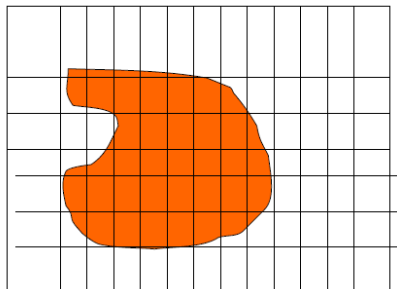
- 1 Markov Chains
- 2 Markov Chain Monte Carlo
- 3 Application to parameter fitting

Applications

(lecture 04)

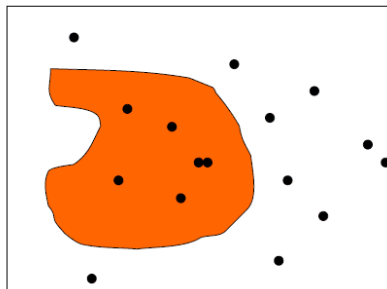
Described procedure can be used not only to calculate integrals of one-dimensional functions, it is much more general...

How to calculate volume of a given shape?



Standard procedure:

scan all dimensions using dense point grid and sum cells with centers inside the volume



Monte Carlo integration:

Generate random points in the considered parameter space and count points inside the volume

General case

Examples presented considered the special case: input random variables had uniform distribution and “test function” was binary (returning 0 or 1).

In the general case we want to determine an expectation value of a function $h(\mathbf{x})$ of random variable vector \mathbf{x} described by $f(\mathbf{x})$ pdf:

$$\mu_h \equiv \mathbb{E}_f[h(\mathbf{x})] = \int d\mathbf{x} h(\mathbf{x}) f(\mathbf{x})$$

Monte Carlo determination of μ_h assumes we can generate random variables according to $f(\mathbf{x})$. We can then calculate:

$$\mu_{MC} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i h(\mathbf{x}_i)$$

where \mathbf{x}_i , $i = 1, \dots, N$ are random (input) variables generated from $f(\mathbf{x})$

Weighted Monte Carlo

General method for generating random points in multi-dimensional space using [acceptance–rejection technique](#) can have very low efficiency, if probability distribution function $f(\mathbf{x})$ varies a lot, eg. has sharp peaks.

Assume we know how to generate random numbers from $g(\mathbf{x})$.

We can then apply the following procedure:

- generate \mathbf{x}_i distributed according to $g(\mathbf{x})$
- accept all generated value \mathbf{x}_i ,
but consider them with additional weight: $w_i = f(\mathbf{x})/g(\mathbf{x})$

For example, when calculating the expectation value of $h(\mathbf{x})$:

$$\mu_{MC} \rightarrow \mu_{wMC} = \frac{\sum_i w_i h(\mathbf{x}_i)}{\sum_i w_i}$$

“unweighted” samples considered previously correspond to $w_i \equiv 1$

Weighted Monte Carlo

When using weighted Monte Carlo “events”, number of events has to be replaced by sum of weights:

$$N \rightarrow N_w = \sum_i w_i$$

Variance of the sum of weights:

$$\mathbb{V}(N_w) = \sum_i w_i^2$$

Statistical power of the weighted Monte Carlo sample is equivalent to:

$$N_{eq} = \frac{N_w^2}{\mathbb{V}(N_w)} = \frac{(\sum_i w_i)^2}{\sum_i w_i^2}$$

For Poisson distributed random number $\mathbb{V}(N) = N$

General problem

Presented above was a simple example of a more general problem: how to estimate parameters of the probability distribution function from the results of the experiment (measurements).

In many cases, parameter value can not be directly extracted from the outcome of the measurement.

In the general case, shape of the probability density function for \mathbf{x} :

$$\mathbf{x} = (x_1, \dots, x_n)$$

depends on a set of pdf parameters:

$$\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$$

so the probability density should be written as:

$$f(\mathbf{x}; \boldsymbol{\lambda})$$

Maximum Likelihood Method

The product:

$$L = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \boldsymbol{\lambda})$$

is called a **likelihood function**.

The most commonly used approach to parameter estimation is the **maximum likelihood approach**:

as the **best estimate of the parameter set $\boldsymbol{\lambda}$** we choose the parameter values for which the **likelihood function has a (global) maximum**.

Frequently used is also log-likelihood function

$$\ell = \ln L = \sum_{j=1}^N \ln f(\mathbf{x}^{(j)}; \boldsymbol{\lambda})$$

we can look for maximum value of ℓ or minimum of $-2\ell = -2\ln L$

Markov Chains

- 1 Markov Chains
- 2 Markov Chain Monte Carlo
- 3 Application to parameter fitting

General concept

(Bonamente)

Markov Chain is a stochastic process where we consider the sequence of measurements (random variables) $X^{(t)}$.

Measurements at fixed time intervals are a frequent case...

Outcome of the measurement (also called “state” of the chain) has to belong to the defined “state space”. It is our sample space...

However, the probability density for different states is not given a priori! Instead, probability of the subsequent state (measurement at $t + 1$) depends only on the current state of the system:

$$P(X^{(t+1)}) = P(X^{(t+1)}|X^{(t)})$$

General concept

(Bonamente)

Markov Chain is a stochastic process where we consider the sequence of measurements (random variables) $X^{(t)}$.

Measurements at fixed time intervals are a frequent case...

Outcome of the measurement (also called “state” of the chain) has to belong to the defined “state space”. It is our sample space...

However, the probability density for different states is not given a priori! Instead, probability of the subsequent state (measurement at $t + 1$) depends only on the current state of the system:

$$P(X^{(t+1)}) = P(X^{(t+1)}|X^{(t)})$$

Probability can change in time, but dependent only on the current state of the chain, and not on any of its previous history!

This “short memory” property is known as the “Markovian property”.

Simple example: Ehrenfest chain

(Bonamente)

Simple model of diffusion: consider two boxes with a total of N balls.

The state of the system can be defined by a number n of balls which are placed in the first box, $0 \leq n \leq N$.

The state space of the system has $N + 1$ elements.

Simple example: Ehrenfest chain

(Bonamente)

Simple model of diffusion: consider two boxes with a total of N balls.

The state of the system can be defined by a number n of balls which are placed in the first box, $0 \leq n \leq N$.

The state space of the system has $N + 1$ elements.

The Ehrenfest chain is defined by the following procedure. At each step:

- select a ball at random from either box,
- place the selected ball in the other box.

Simple example: Ehrenfest chain

(Bonamente)

Simple model of diffusion: consider two boxes with a total of N balls.

The state of the system can be defined by a number n of balls which are placed in the first box, $0 \leq n \leq N$.

The state space of the system has $N + 1$ elements.

The Ehrenfest chain is defined by the following procedure. At each step:

- select a ball at random from either box,
- place the selected ball in the other box.

This chain can be presented in terms of the transition probabilities:

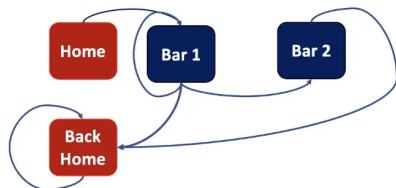
$$p(n^{(t+1)}) = \begin{cases} \frac{n^{(t)}}{N} & \text{for } n^{(t+1)} = n^{(t)} - 1 \\ \frac{N-n^{(t)}}{N} & n^{(t+1)} = n^{(t)} + 1 \\ 0 & n^{(t+1)} \neq n^{(t)} \pm 1 \end{cases}$$

Web example

Piero Paialunga in Towards Data Science

As a student you can go to the bar each Saturday.

And you need to go back home at some time...



We can consider the following “chain” of states (shown above):

- you always start from Home going to Bar 1 or Bar 2.
- after each drink in Bar 1 you have three choices: go Back Home, go to Bar 2 and order another drink in Bar 1.
- if you are already in Bar 2, you have only two choices after each round: go Back Home or order another drink (not shown).
- once you get Back Home, you stay there.

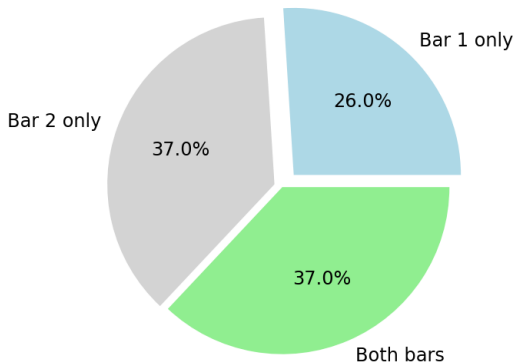
Web example

Piero Paialunga in Towards Data Science

Even if all transition probabilities are known, it is not always possible to obtain statistical properties of the distribution directly...

But one can simulate Markov Chain state sequence many times...

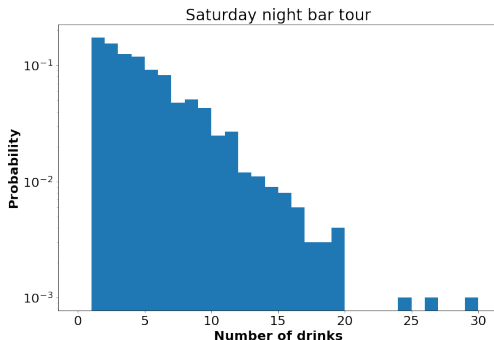
Probability of visiting bars:



Web example

Piero Paialunga in Towards Data Science

Probability density for the number of drinks:

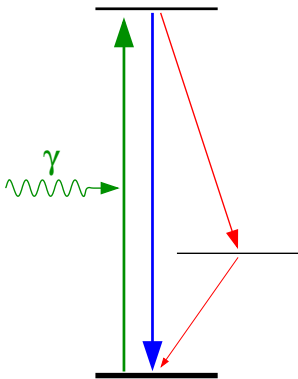


We can not only estimate the expected number of drinks (which we could also do from the known probabilities), but also the distribution...

Another example

The chain in the web example always ended in the single 'Back Home' state.

Not very interesting...

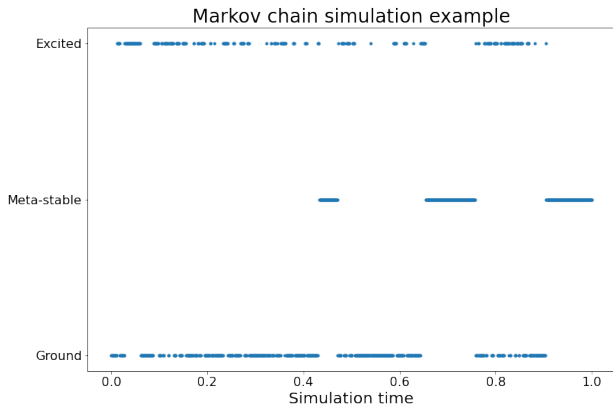


Consider an atom irradiated with laser light tuned to the excitation energy:

- when in ground state, atom has certain probability (per time unit) to get excited
- when in the excited state, atom can radiate photon and go back to the ground state or, with lower probability, radiate softer photon and go to intermediate meta-stable state.
- when in the meta-stable state, probability of radiation (per unit of time) is very low.

Another example

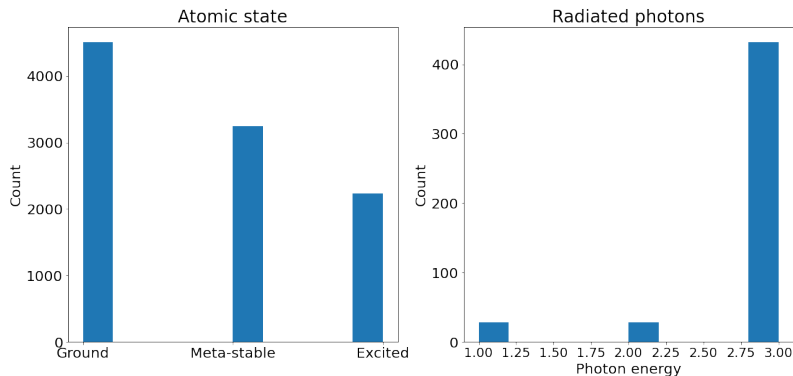
Example simulation results starting from ground state, 1000 time steps:



Fast oscillations between ground and excited state, longer stays in meta-stable...

Another example

Example simulation results starting from ground state, 10000 time steps:

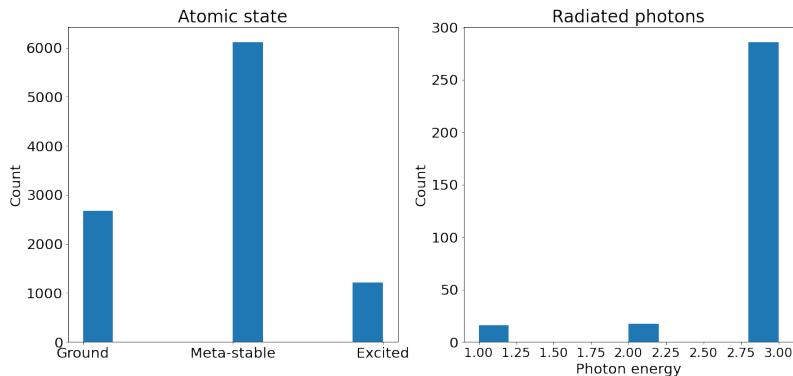


System "forgets" about the initial state after few time steps.

We can get distributions for different parameters...

Another example

Example simulation results starting from ground state, 10000 time steps:
After increasing meta-stable state lifetime:



System "forgets" about the initial state after few time steps.

We can get distributions for different parameters...

Transition probability

Assume that the state space consists of N states: $s_{(1)}, \dots, s_{(N)}$. Then, for each state $s_{(i)}$ one can define a set of on-step transition probabilities:

$$p_{ij} = p(X^{(t+1)} = s_{(j)} | X^{(t)} = s_{(i)})$$

We usually require that these probabilities are time-independent (such chain is called time-homogeneous).

Transition probability

Assume that the state space consists of N states: $s_{(1)}, \dots, s_{(N)}$. Then, for each state $s_{(i)}$ one can define a set of on-step transition probabilities:

$$p_{ij} = p(X^{(t+1)} = s_{(j)} | X^{(t)} = s_{(i)})$$

We usually require that these probabilities are time-independent (such chain is called time-homogeneous).

If we now describe state of the system by a N -component vector:

$$(s_{(i)})_j = \delta_{ij} \quad \text{e.g. } s_{(1)} = (1, 0, 0, \dots, 0)$$

Transition probability

Assume that the state space consists of N states: $s_{(1)}, \dots, s_{(N)}$. Then, for each state $s_{(i)}$ one can define a set of on-step transition probabilities:

$$p_{ij} = p(X^{(t+1)} = s_{(j)} | X^{(t)} = s_{(i)})$$

We usually require that these probabilities are time-independent (such chain is called time-homogeneous).

If we now describe state of the system by a N -component vector:

$$(s_{(i)})_j = \delta_{ij} \quad \text{e.g. } s_{(2)} = (0, 1, 0, \dots, 0)$$

Transition probability

Assume that the state space consists of N states: $s_{(1)}, \dots, s_{(N)}$. Then, for each state $s_{(i)}$ one can define a set of on-step transition probabilities:

$$p_{ij} = p(X^{(t+1)} = s_{(j)} | X^{(t)} = s_{(i)})$$

We usually require that these probabilities are time-independent (such chain is called time-homogeneous).

If we now describe state of the system by a N -component vector:

$$(s_{(i)})_j = \delta_{ij} \quad \text{e.g. } s_{(N)} = (0, 0, 0, \dots, 1)$$

Transition probability

Assume that the state space consists of N states: $s_{(1)}, \dots, s_{(N)}$. Then, for each state $s_{(i)}$ one can define a set of on-step transition probabilities:

$$p_{ij} = p(X^{(t+1)} = s_{(j)} | X^{(t)} = s_{(i)})$$

We usually require that these probabilities are time-independent (such chain is called time-homogeneous).

If we now describe state of the system by a N -component vector:

$$(s_{(i)})_j = \delta_{ij} \quad \text{e.g. } s_{(N)} = (0, 0, 0, \dots, 1)$$

then probabilities for different states to proceed after state $s_{(i)}$ can be written as:

$$p = s_{(i)} \cdot \mathbb{T} \quad \text{where } \mathbb{T} = (p_{ij})$$

is the transition matrix

Chain properties

(Bonamente)

Probabilities of states after n time steps are then given by:

$$p^{(n)} = s_{(i)} \cdot \mathbb{T}^n$$

Chain properties

(Bonamente)

Probabilities of states after n time steps are then given by:

$$p^{(n)} = s_{(i)} \cdot \mathbb{T}^n$$

Let u_k denote the probability that the system returns to the initial state $s_{(i)}$ in exactly k time steps. We can define the total probability for returning to the initial state:

$$u = \sum_{k=1}^{\infty} u_k$$

Chain properties

(Bonamente)

Probabilities of states after n time steps are then given by:

$$p^{(n)} = s_{(i)} \cdot \mathbb{T}^n$$

Let u_k denote the probability that the system returns to the initial state $s_{(i)}$ in exactly k time steps. We can define the total probability for returning to the initial state:

$$u = \sum_{k=1}^{\infty} u_k$$

States can be classified according to this probability:

- if $u = 1$ state $s_{(i)}$ is recurrent,
- if $u < 1$ state $s_{(i)}$ is transient.

If state is recurrent, it will certainly be observed again (even, if we have to wait very long), and the system will return to this state infinitely often.

Chain properties

(Bonamente)

State $s_{(j)}$ is **accessible** from the initial state $s_{(i)}$, if there is a non-zero probability of reaching this state from the initial state in finite number of time steps:

$$\left(p^{(m)}\right)_j = \left(s_{(i)} \cdot \mathbb{T}^m\right)_j > 0$$

for some natural number m .

Chain properties

(Bonamente)

State $s_{(j)}$ is **accessible** from the initial state $s_{(i)}$, if there is a non-zero probability of reaching this state from the initial state in finite number of time steps:

$$\left(p^{(m)}\right)_j = \left(s_{(i)} \cdot \mathbb{T}^m\right)_j > 0$$

for some natural number m .

If a state $s_{(j)}$ is accessible from a recurrent state $s_{(i)}$, then $s_{(j)}$ is also recurrent, and $s_{(i)}$ is accessible from $s_{(j)}$.

Chain properties

(Bonamente)

State $s_{(j)}$ is **accessible** from the initial state $s_{(i)}$, if there is a non-zero probability of reaching this state from the initial state in finite number of time steps:

$$\left(p^{(m)}\right)_j = \left(s_{(i)} \cdot \mathbb{T}^m\right)_j > 0$$

for some natural number m .

If a state $s_{(j)}$ is accessible from a recurrent state $s_{(i)}$, then $s_{(j)}$ is also recurrent, and $s_{(i)}$ is accessible from $s_{(j)}$.

If a Markov chain has a finite number of states and each state is accessible from any other state, then all states are recurrent.

Chain properties

(Bonamente)

A chain is said to be **irreducible** if all states are accessible from others.

Possible states of reducible Markov Chain can be divided into two or more classes, which do not communicate with each other.

Chain properties

(Bonamente)

A chain is said to be **irreducible** if all states are accessible from others. Possible states of reducible Markov Chain can be divided into two or more classes, which do not communicate with each other.

A state $s_{(i)}$ is said to be **periodic with period T** if system can return to this state only at times t divisible by T :

$$\left(p^{(t)}\right)_j = \begin{cases} p > 0 & \text{for } t \% T = 0 \\ 0 & \text{for } t \% T \neq 0 \end{cases}$$

All states of irreducible chain share the same period.

Chain properties

(Bonamente)

A chain is said to be **irreducible** if all states are accessible from others. Possible states of reducible Markov Chain can be divided into two or more classes, which do not communicate with each other.

A state $s_{(i)}$ is said to be **periodic with period T** if system can return to this state only at times t divisible by T :

$$\left(p^{(t)}\right)_j = \begin{cases} p > 0 & \text{for } t \% T = 0 \\ 0 & \text{for } t \% T \neq 0 \end{cases}$$

All states of irreducible chain share the same period.

A chain is said to be **aperiodic** if return to a given state can occur at any time (corresponding to $T = 1$ in definition above).

Stationary distribution

In most cases, we do not care about the initial system state, we want to calculate the set of probabilities for a system after a large number n of steps:

$$p^\infty = \lim_{n \rightarrow \infty} p^{(n)}$$

This probabilities are called **limiting probabilities**.

Stationary distribution

In most cases, we do not care about the initial system state, we want to calculate the set of probabilities for a system after a large number n of steps:

$$p^\infty = \lim_{n \rightarrow \infty} p^{(n)}$$

This probabilities are called **limiting probabilities**.

For a **irreducible aperiodic** Markov Chain with **recurrent** states, limiting probabilities correspond to the **stationary distribution**:

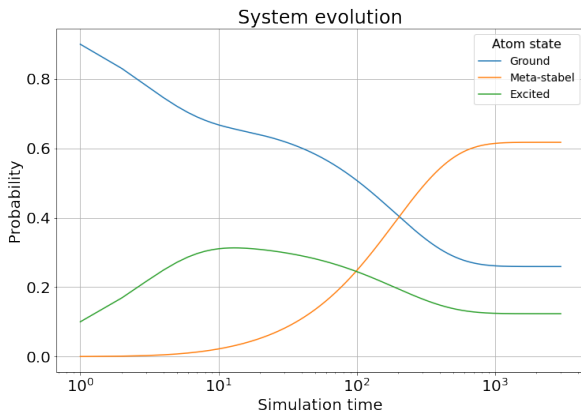
$$\pi = \pi \cdot \mathbb{T}$$

and that this distribution is **unique**.

Regardless of the starting point of the chain, the same stationary distribution will eventually be reached.

Stationary distribution

Evolution of state probabilities for system starting at 'Ground' state at $t = 0$



Stationary state reached for $t \sim 1000$

Note logarithmic time scale!

Stationary distribution

this is what we look for in most cases

There are three possible approaches to finding a stationary solution:

- by running multiple Markov Chain instances and looking at final state distribution, **simple but time consuming**

Stationary distribution

this is what we look for in most cases

There are three possible approaches to finding a stationary solution:

- by running multiple Markov Chain instances and looking at final state distribution, **simple but time consuming**
- by taking arbitrary initial state probability vector and applying the transfer matrix many times,

Stationary distribution

this is what we look for in most cases

There are three possible approaches to finding a stationary solution:

- by running multiple Markov Chain instances and looking at final state distribution, **simple but time consuming**
- by taking arbitrary initial state probability vector and applying the transfer matrix many times,
- by looking for analytic solution to the problem:

$$\pi_j = \sum_j \pi_j p_{ij} \quad \text{stationary distribution}$$

$$\sum_i \pi_i = 1 \quad \text{normalization constrain}$$

$$\pi_j \geq 0$$

Stationary distribution

Herman Scheepers on Towards Data Science

In the analytic approach the problem can be presented as a set of equations:

$$\begin{pmatrix} \mathbb{T}^T - \mathbb{I} \\ \hline 1 \quad \dots \quad 1 \end{pmatrix} \cdot \boldsymbol{\pi} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \hline 1 \end{pmatrix}$$

A · π = b

which are, however, not independent (the problem is over-constrained).

Stationary distribution

Herman Scheepers on Towards Data Science

In the analytic approach the problem can be presented as a set of equations:

$$\begin{pmatrix} \mathbb{T}^\top - \mathbb{I} \\ \hline 1 \quad \dots \quad 1 \end{pmatrix} \cdot \boldsymbol{\pi} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \hline 1 \end{pmatrix}$$

$$\mathbb{A} \cdot \boldsymbol{\pi} = \mathbf{b}$$

which are, however, not independent (the problem is over-constrained).

The simple solution is to multiply both sides by \mathbb{A}^\top :

$$\mathbb{A}^\top \mathbb{A} \cdot \boldsymbol{\pi} = \mathbb{A}^\top \mathbf{b}$$

which can now be solved with standard linear algebra procedures...

Markov Chains

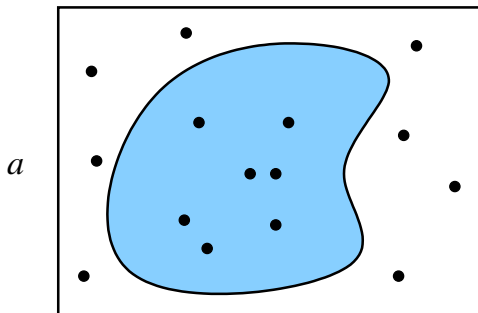
- 1 Markov Chains
- 2 Markov Chain Monte Carlo
- 3 Application to parameter fitting

General concept

(Katzgraber, arXiv:0905.1629)

We introduced Monte Carlo as an alternative method for integrating an arbitrary function.

Arbitrary parameter space can be considered.



Rejection technique

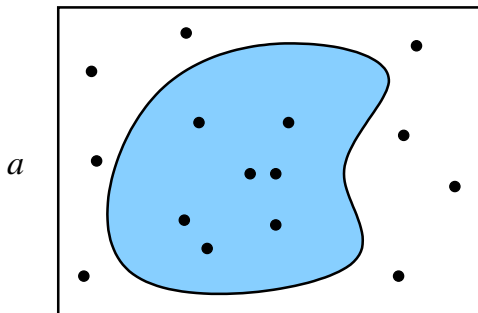
Generate uniformly distributed random points and select those in the considered parameter space...

General concept

(Katzgraber, arXiv:0905.1629)

We introduced Monte Carlo as an alternative method for integrating an arbitrary function.

Arbitrary parameter space can be considered.



a

b

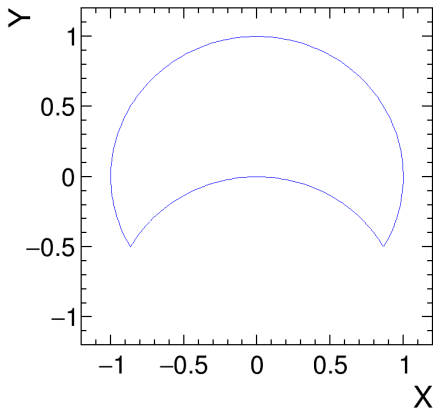
Rejection technique

Generate uniformly distributed random points and select those in the considered parameter space...

Efficiency can be low...

Standard approach example

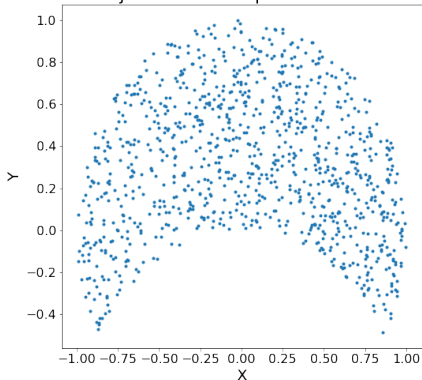
Generation of random points from the surface considered in lecture 04



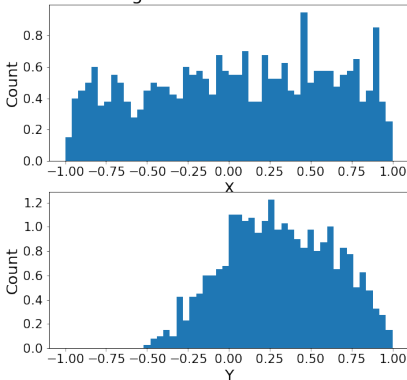
Standard approach example

Generation of random points from the surface considered in lecture 04

Rejection technique simulation



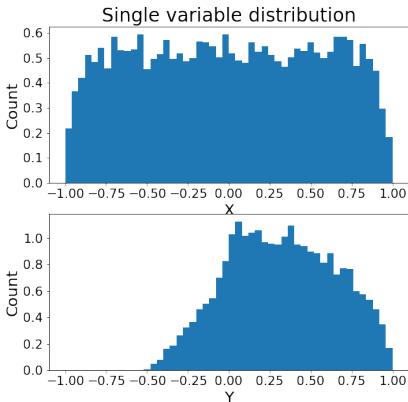
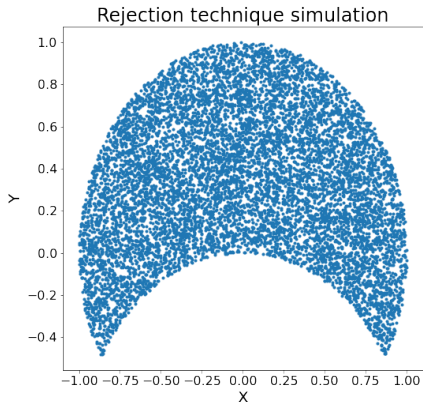
Single variable distribution



N=1 000

Standard approach example

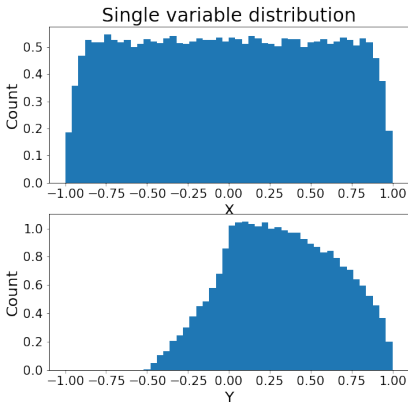
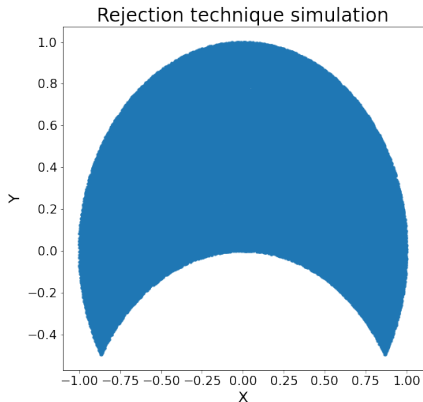
Generation of random points from the surface considered in lecture 04



N=10 000

Standard approach example

Generation of random points from the surface considered in lecture 04

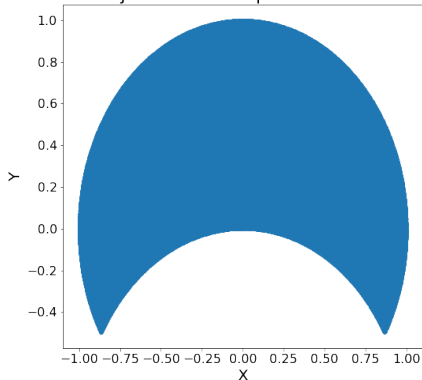


$N=100\ 000$

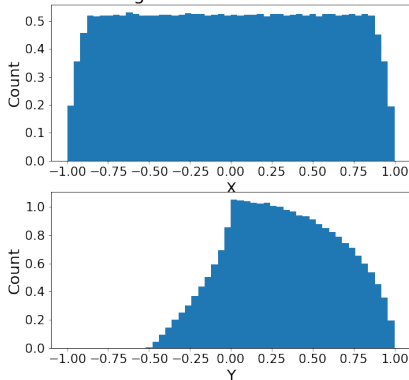
Standard approach example

Generation of random points from the surface considered in lecture 04

Rejection technique simulation



Single variable distribution



$N=1\ 000\ 000$

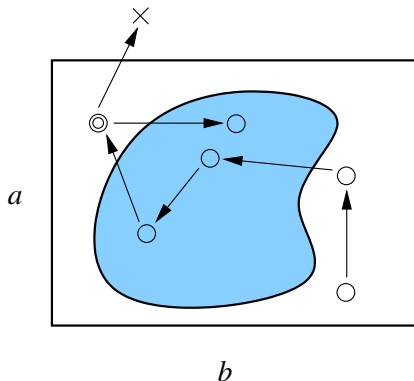
generated in 2 093 551 tries

General concept

(Katzgraber, arXiv:0905.1629)

We do not want to reject events!

Random move procedure:
subsequent points generated by
random variations of previous ones



Markov Chain Monte Carlo procedure

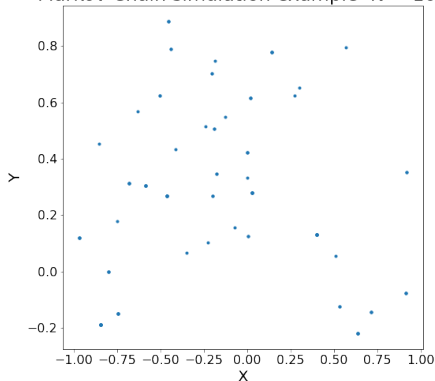
If the new point is outside the considered parameter space,
do not reject it, but **take the last point again** (!)

Can this procedure work ?

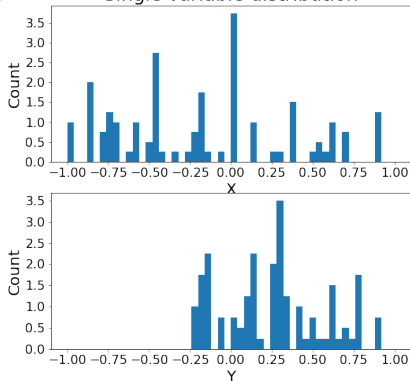
Markov Chain MC example

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example N = 100



Single variable distribution

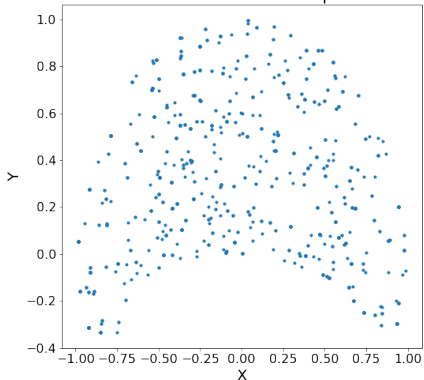


N=100

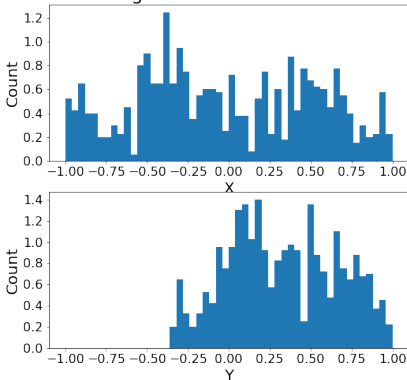
Markov Chain MC example

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example N = 1000



Single variable distribution

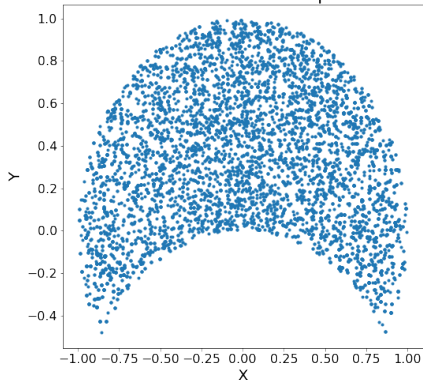


N=1 000 Fluctuations are larger, as many points “duplicated”

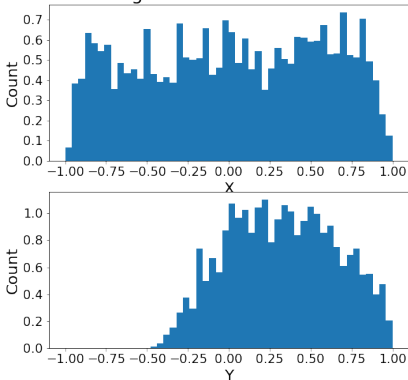
Markov Chain MC example

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example N = 10000



Single variable distribution

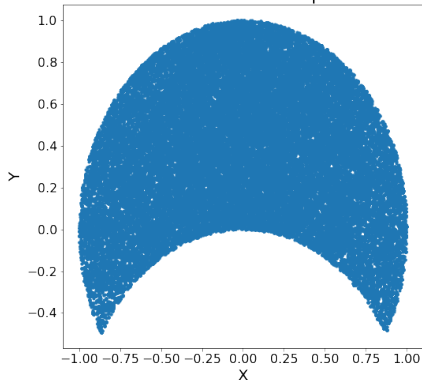


N=10 000

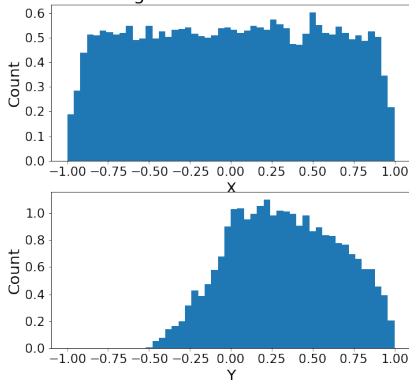
Markov Chain MC example

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example $N = 100000$



Single variable distribution

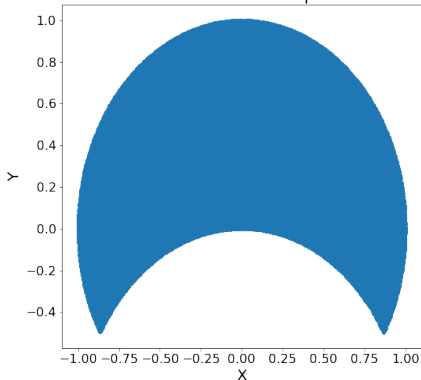


$N=100\ 000$

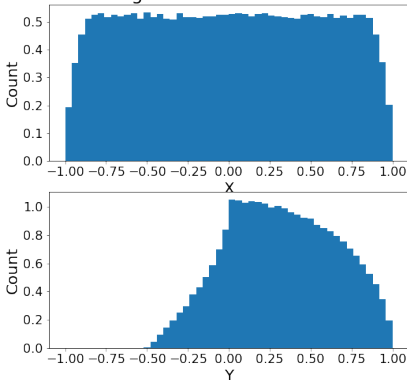
Markov Chain MC example

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example $N = 1000000$



Single variable distribution



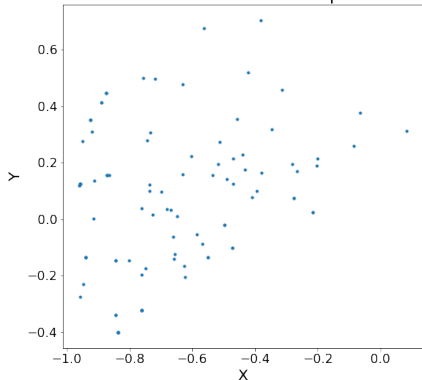
$N=1\,000\,000$

But “duplicates” not relevant for $N \rightarrow \infty$

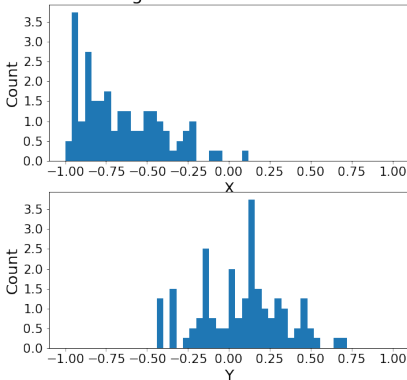
Markov Chain example

We can reduce number of “duplicates” by reducing step: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 100



Single variable distribution

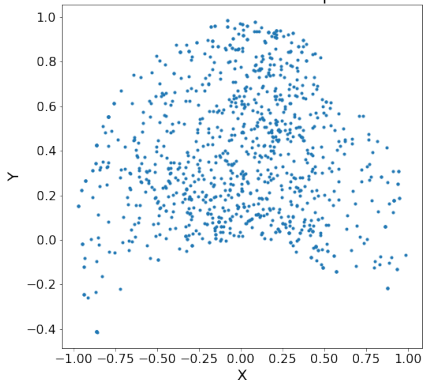


N=100 Significant bias, depending on the starting point...

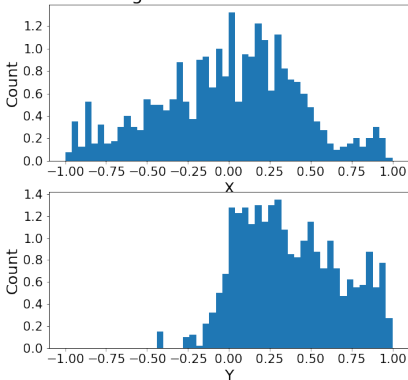
Markov Chain example

We can reduce number of “duplicates” by reducing step: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 1000



Single variable distribution

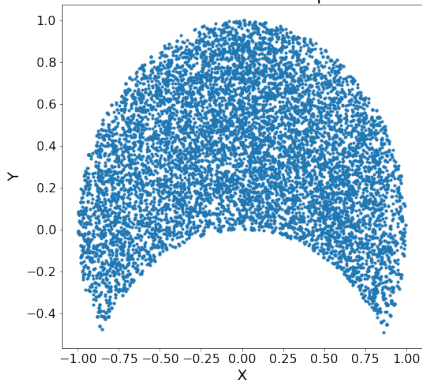


N=1 000

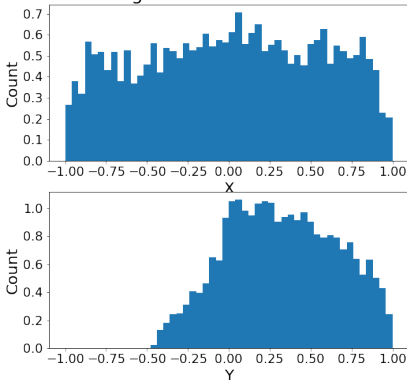
Markov Chain example

We can reduce number of “duplicates” by reducing step: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 10000



Single variable distribution

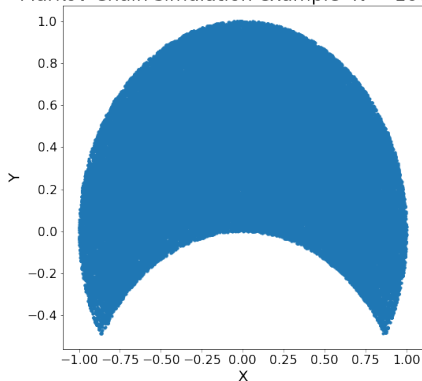


N=10 000

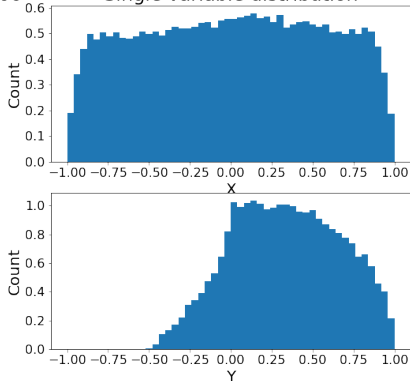
Markov Chain example

We can reduce number of “duplicates” by reducing step: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 100000



Single variable distribution



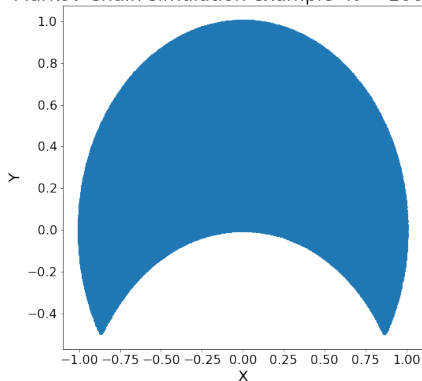
N=100 000

Distribution still not uniform...

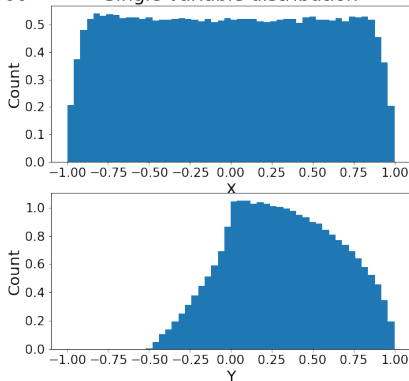
Markov Chain example

We can reduce number of “duplicates” by reducing step: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example $N = 1000000$



Single variable distribution



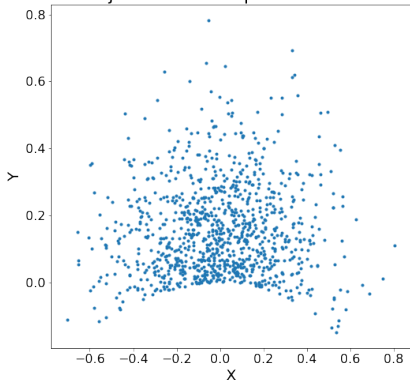
$N=1\ 000\ 000$

But gets uniform for $N \rightarrow \infty$

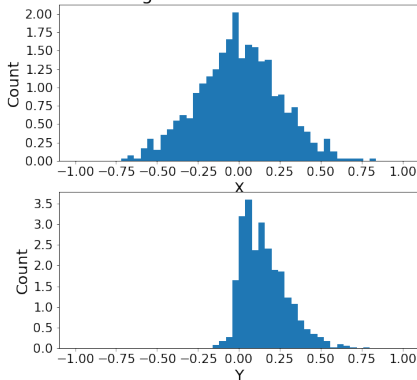
More general case

Gaussian probability distribution in the considered parameter space

Rejection technique simulation



Single variable distribution

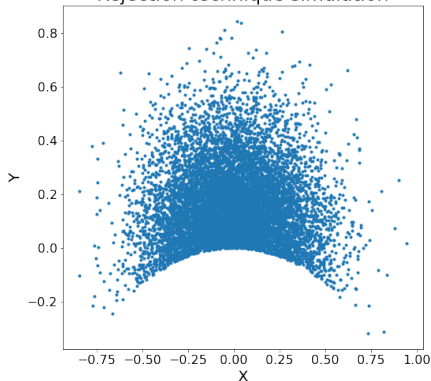


N=1 000

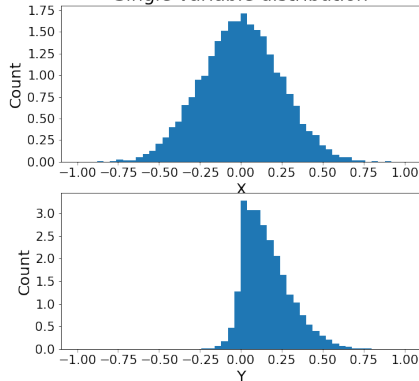
More general case

Gaussian probability distribution in the considered parameter space

Rejection technique simulation



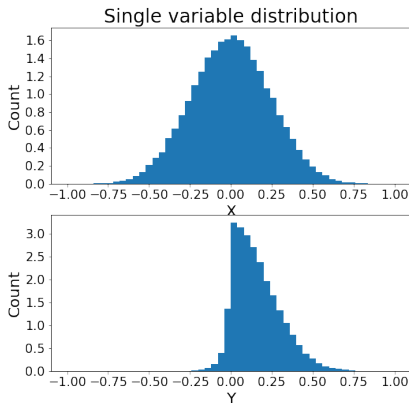
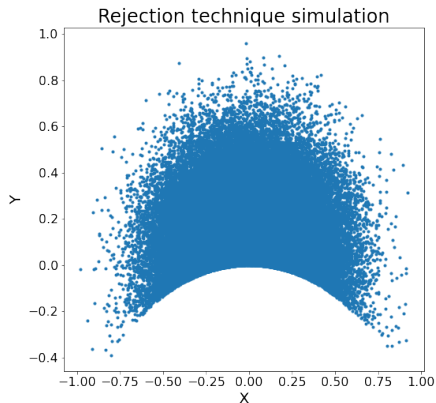
Single variable distribution



N=10 000

More general case

Gaussian probability distribution in the considered parameter space



N=100 000 generated in 2 335 937 tries, 4.3% efficiency

Metropolis–Hastings algorithm

(Givens)

Consider chain described by on-step transition probability $p(X^{(t+1)}|X^{(t)})$

To generate points distributed according to $f(X)$, for each step t :

- generate candidate point X^* from $p(X^*|X^{(t)})$
- compute the Metropolis–Hastings ratio:

$$R = \frac{f(X^*) p(X^{(t)}|X^*)}{f(X^{(t)}) p(X^*|X^{(t)})}$$

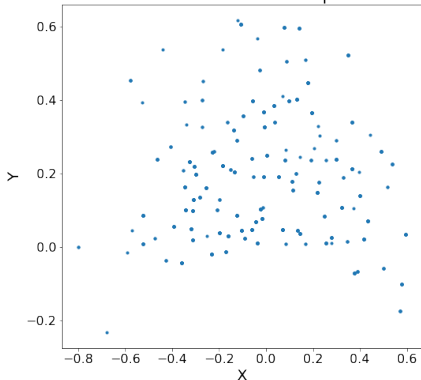
- for the next step take

$$X^{(t+1)} = \begin{cases} X^* & \text{with probability } p^* = \min\{R, 1\} \\ X^{(t)} & \text{otherwise.} \end{cases} \quad \text{with probability } 1 - p^*$$

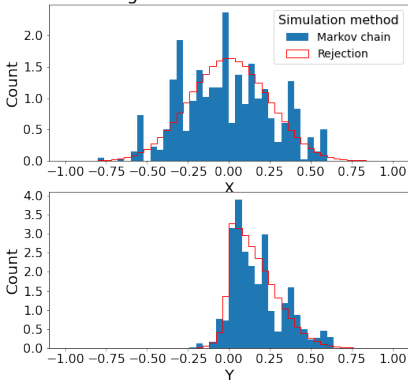
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example N = 1000



Single variable distributions



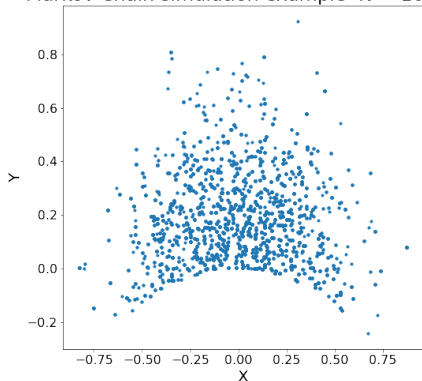
N=1 000

Large step \Rightarrow large fluctuations

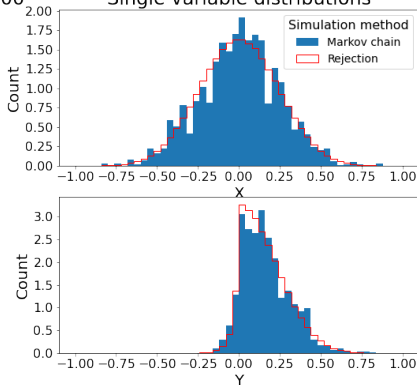
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example N = 10000



Single variable distributions



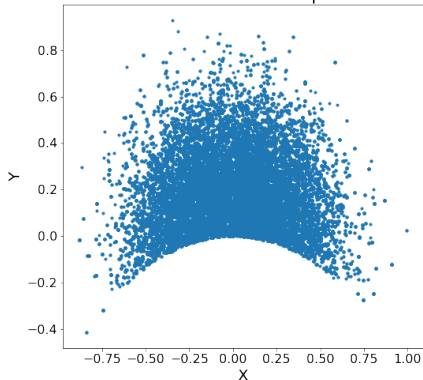
N=10 000

Large step \Rightarrow large fluctuations

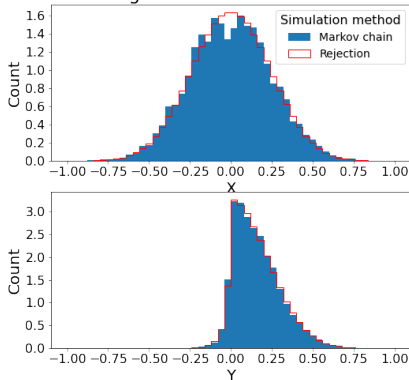
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example $N = 100000$



Single variable distributions



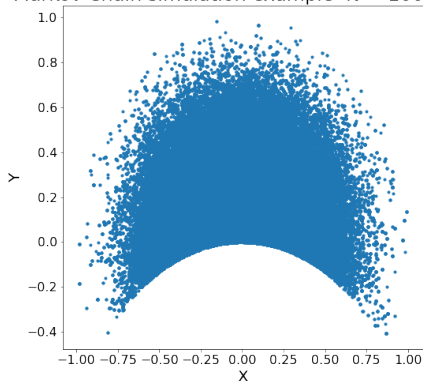
$N=100\ 000$

But converges to the expected distribution for large N

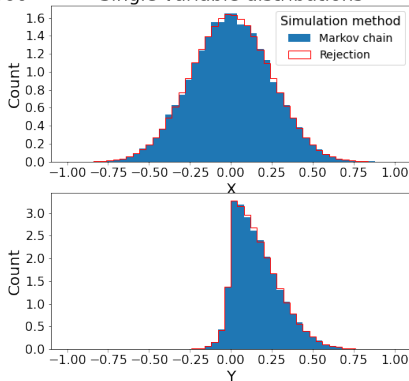
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 1$

Markov Chain simulation example $N = 1000000$



Single variable distributions



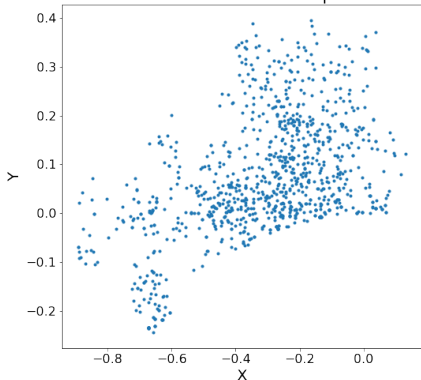
$N=1\ 000\ 000$

But converges to the expected distribution for large N

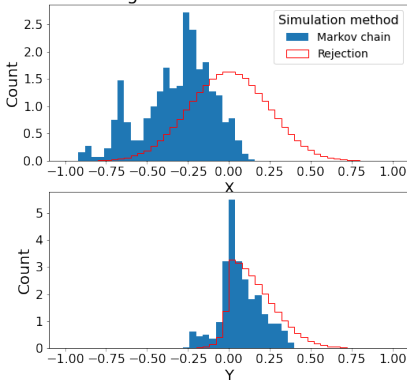
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.05$

Markov Chain simulation example N = 1000



Single variable distributions



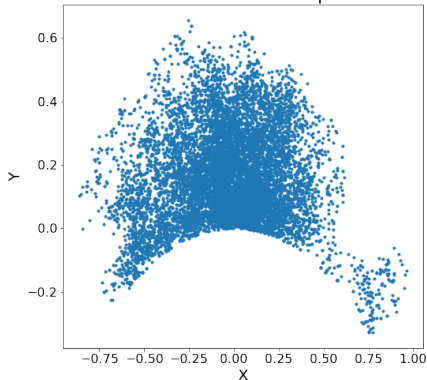
N=1 000

Small step \Rightarrow large bias

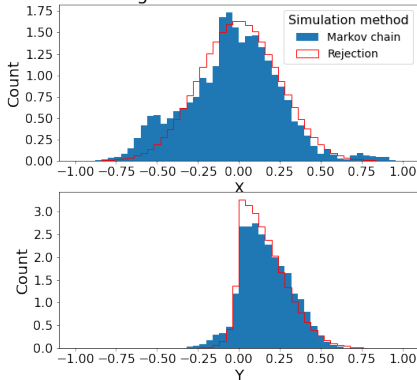
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.05$

Markov Chain simulation example N = 10000



Single variable distributions



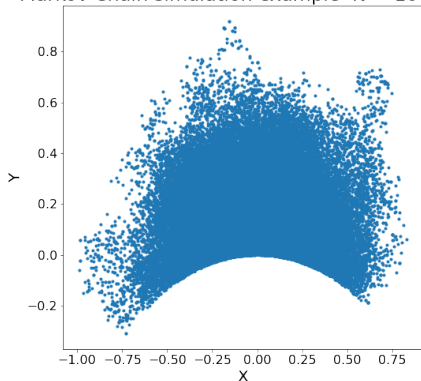
N=10 000

Small step \Rightarrow large bias

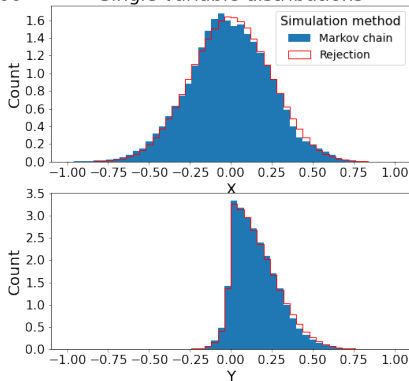
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.05$

Markov Chain simulation example $N = 100000$



Single variable distributions



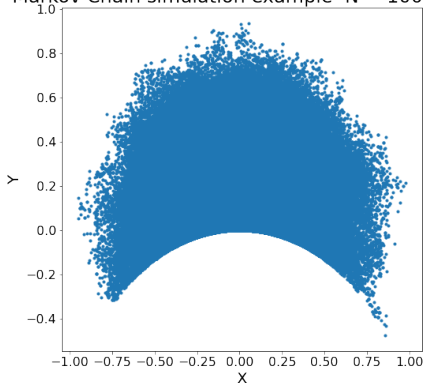
$N=100\ 000$

But converges to the expected distribution for large N

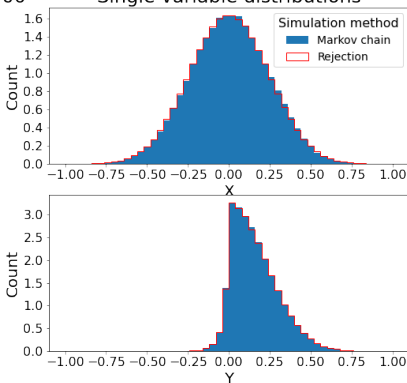
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.05$

Markov Chain simulation example $N = 1000000$



Single variable distributions



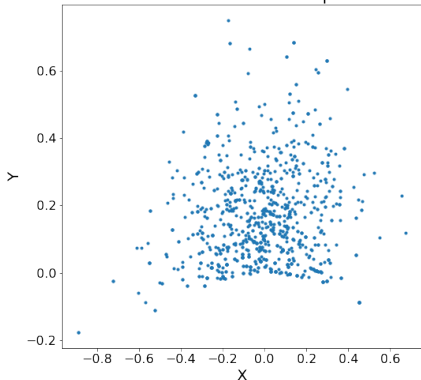
$N=1\ 000\ 000$

But converges to the expected distribution for large N

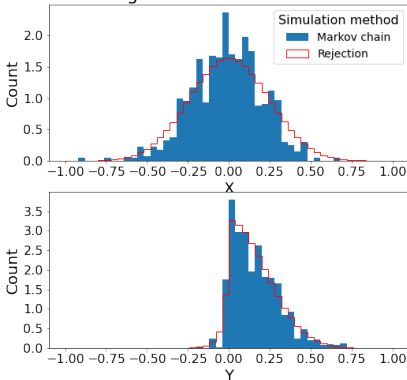
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 1000



Single variable distributions



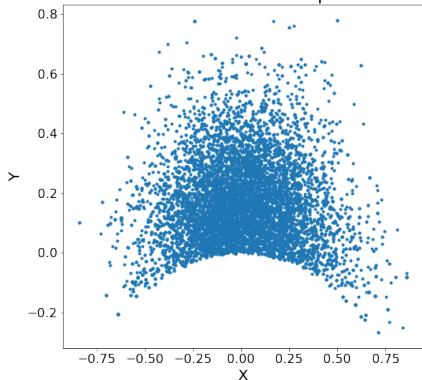
N=1 000

Optimal step $\Rightarrow \sim$ Poisson fluctuations, minimum bias

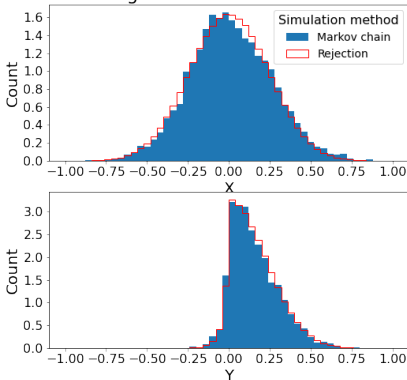
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 10000



Single variable distributions



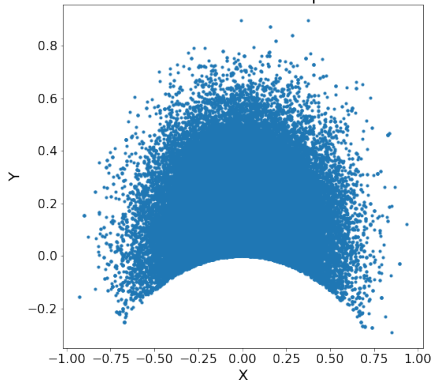
N=10 000

Optimal step $\Rightarrow \sim$ Poisson fluctuations, minimum bias

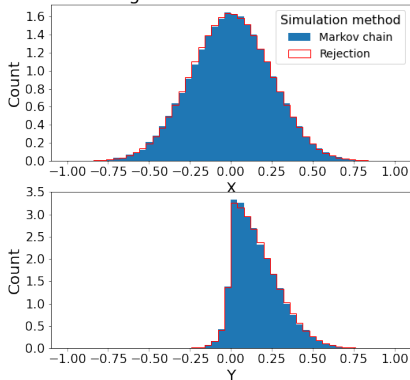
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example N = 100000



Single variable distributions



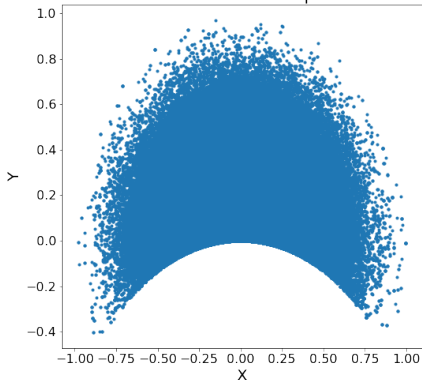
N=100 000

Converges fast to the expected distribution

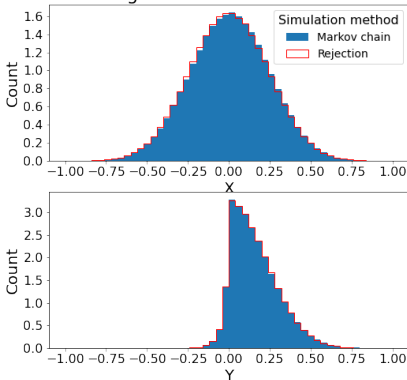
Markov Chain MC example (2)

Using maximum step size: $\Delta x = \Delta y = 0.2$

Markov Chain simulation example $N = 1000000$



Single variable distributions



$N=1\ 000\ 000$

No rejection! Much larger samples with the same CPU

Markov Chains

- 1 Markov Chains
- 2 Markov Chain Monte Carlo
- 3 Application to parameter fitting

Bayesian approach

(lecture 01)

Bayes theorem can be used to generalize the concept of probability. In particular, one can consider “probability” of given hypothesis H (theoretical model or model parameter, eg. Hubble constant) when taking into known outcome D (data) of the experiment

$$P(H|D) = \frac{P(D|H)}{P(D)} \cdot P(H)$$

There are two problems with this approach:

- H can not be considered an event, sampling space can not be defined (no experiment to repeat)
- we need to make a **subjective** assumption about the “prior” $P(H)$ describing our initial belief in hypothesis H

For these reasons I rather use term “degree of belief” for the result of the Bayesian procedure applied to non random events

Bayesian approach

The likelihood function:

$$L(\mathbf{x}, \lambda) = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \lambda)$$

describes the probability of given set of measurement results \mathbf{x} .

Bayesian approach

The likelihood function:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \boldsymbol{\lambda})$$

describes the probability of given set of measurement results \mathbf{x} .

However, in the **bayesian approach** we can use it to construct “probability distribution” for the model parameters $\boldsymbol{\lambda}$:

$$f(\boldsymbol{\lambda}) \sim L(\mathbf{x}, \boldsymbol{\lambda}) \cdot p(\boldsymbol{\lambda})$$

where $p(\boldsymbol{\lambda})$ is the prior distribution for parameters $\boldsymbol{\lambda}$.

If we know $f(\boldsymbol{\lambda})$, we can construct Markov Chain in $\boldsymbol{\lambda}$ space.

Bayesian approach

The likelihood function:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \boldsymbol{\lambda})$$

describes the probability of given set of measurement results \mathbf{x} .

However, in the **bayesian approach** we can use it to construct “probability distribution” for the model parameters $\boldsymbol{\lambda}$:

$$f(\boldsymbol{\lambda}) \sim L(\mathbf{x}, \boldsymbol{\lambda}) \cdot p(\boldsymbol{\lambda})$$

where $p(\boldsymbol{\lambda})$ is the prior distribution for parameters $\boldsymbol{\lambda}$.

If we know $f(\boldsymbol{\lambda})$, we can construct Markov Chain in $\boldsymbol{\lambda}$ space.

With Metropolis–Hastings algorithm, starting from arbitrary $\boldsymbol{\lambda}^{(0)}$ point, the chain should converge to $f(\boldsymbol{\lambda})$ distribution for $N \rightarrow \infty$.

Example

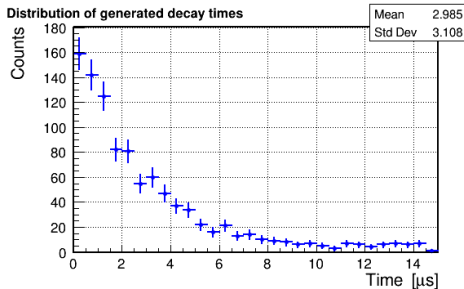
(Homework 10)

1000 events were collected in the **muon lifetime measurement**.

Distribution can be described by the formula:

$$N(t) = \frac{N_{sig}}{\tau} e^{-\frac{t}{\tau}} + N_{bg}$$

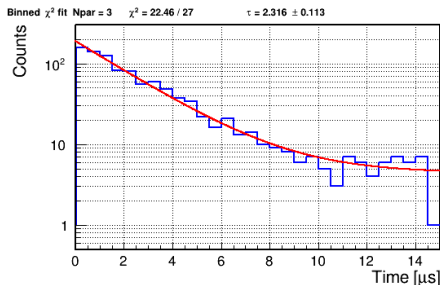
with **flat background** level known to be $N_{bg} = 5 \pm \Delta/2$ (for $\Delta t = 0.5 \mu s$)



Example

(Homework 10)

Histogram can be fitted using **iterative χ^2 minimization** procedure

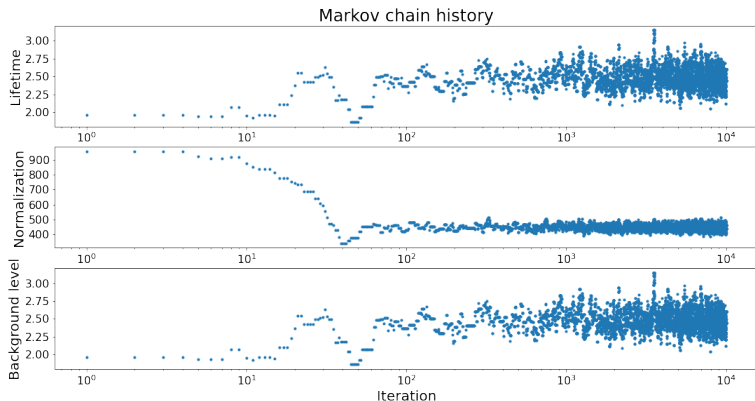


Fit results:

$$\begin{aligned}
 \tau &= 2.316 \pm 0.113 \mu\text{s} \\
 N_{sig} &= 430.773 \pm 16.611 \\
 N_{bg} &= 4.399 \pm 0.424 \\
 \chi^2 &= 22.460/27
 \end{aligned}
 \quad
 \text{Corr} =
 \begin{pmatrix}
 1. & 0.279 & -0.392 \\
 0.279 & 1. & -0.309 \\
 -0.392 & -0.309 & 1.
 \end{pmatrix}$$

Example

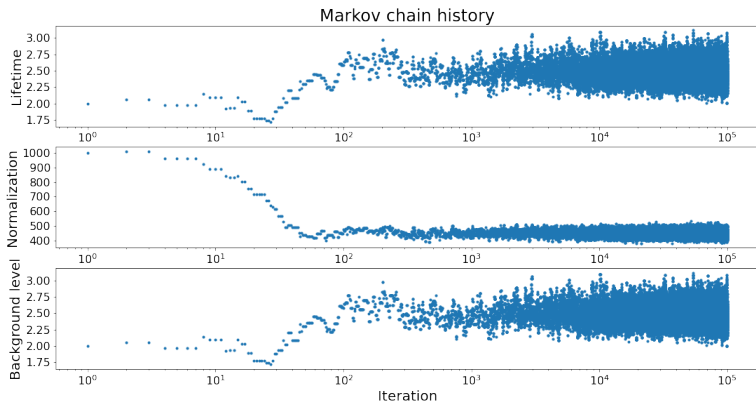
Parameter evolution in the Markov Chain



Stable distribution obtained already after about 100 iterations

Example

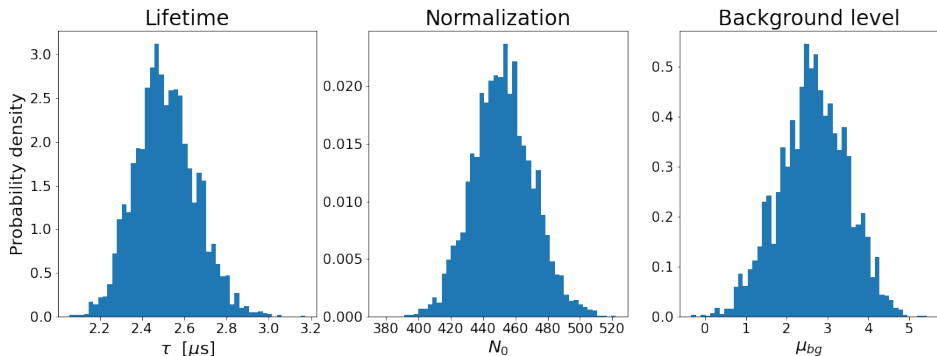
Parameter evolution in the Markov Chain



Stable distribution obtained already after about 100 iterations

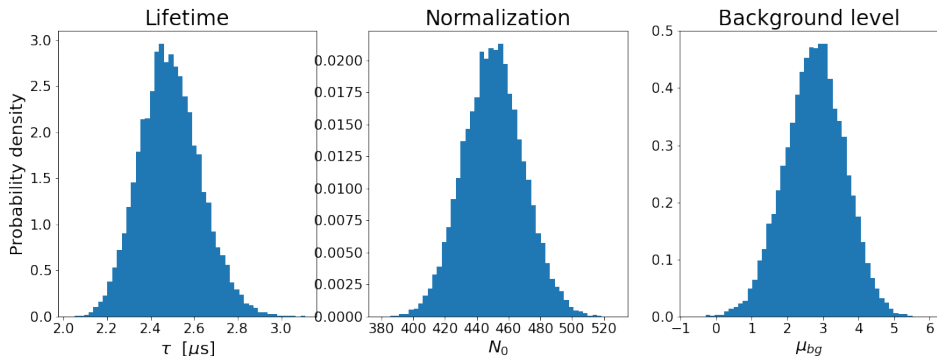
Example

Parameter distributions after $N = 10\,000$ iterations (skipping first 100)



Example

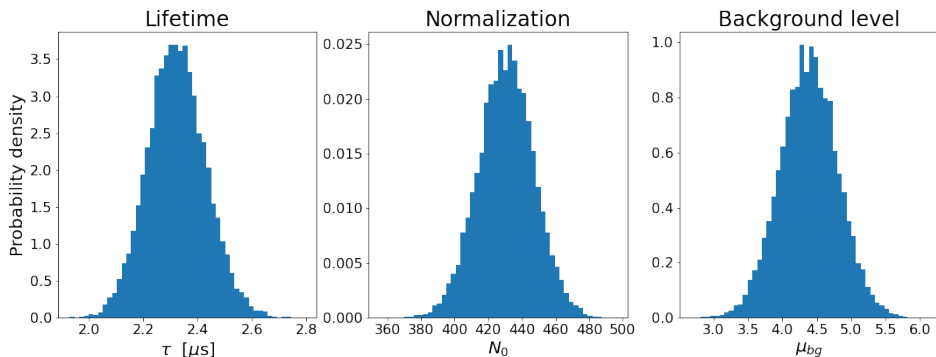
Parameter distributions after $N = 100\,000$ iterations (skipping first 1000)



We can extract expected parameter values with uncertainties...
but also identify problems, e.g. find multiple solutions...

Example

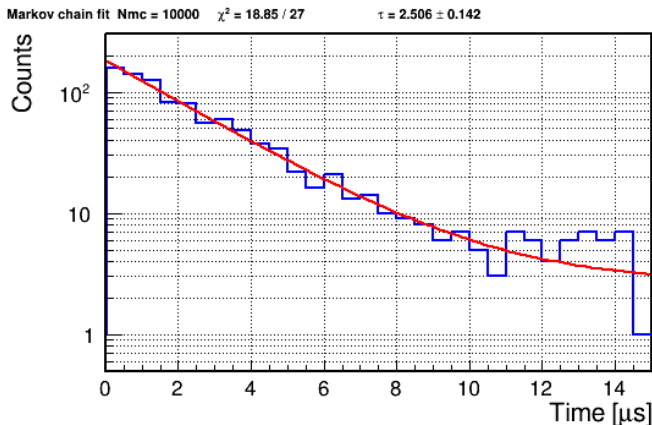
Parameter distributions after $N = 100\,000$ iterations (skipping first 1000)
Including background level constraint



We can extract expected parameter values with uncertainties...
but also identify problems, e.g. find multiple solutions...

Example

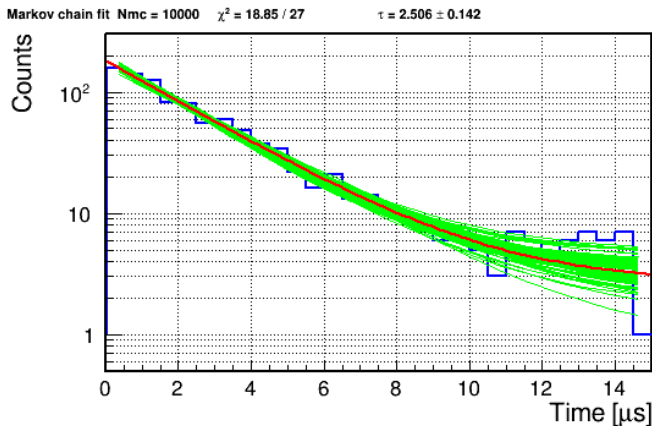
Nominal solution from Markov Chain (mean values of parameters)



Without background constraint

Example

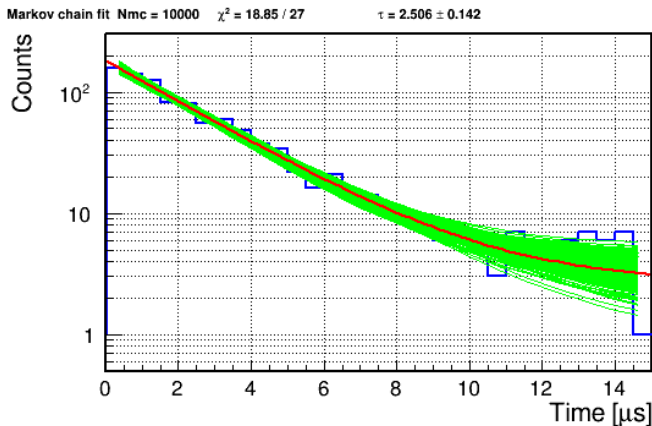
But we can also get the probability distribution of the fit results:



Last 100 chain elements

Example

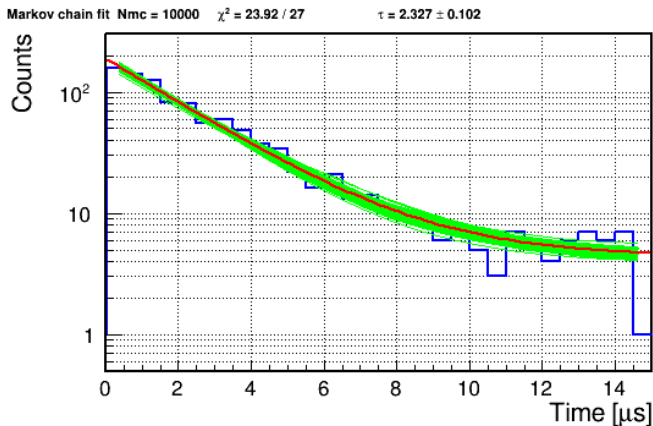
But we can also get the probability distribution of the fit results:



Last 1000 chain elements

Example

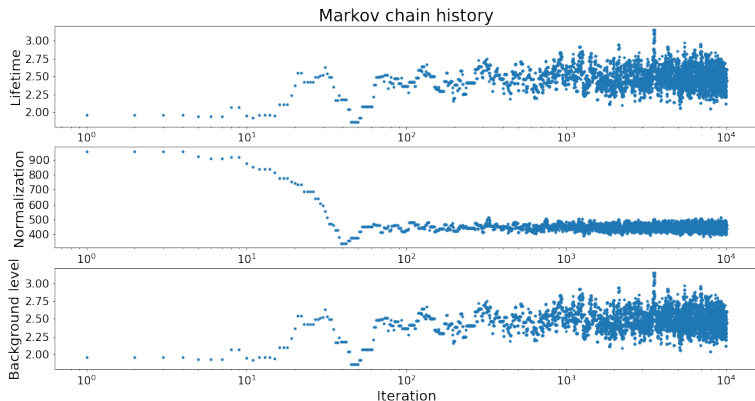
But we can also get the probability distribution of the fit results:



After adding background constraint

Example

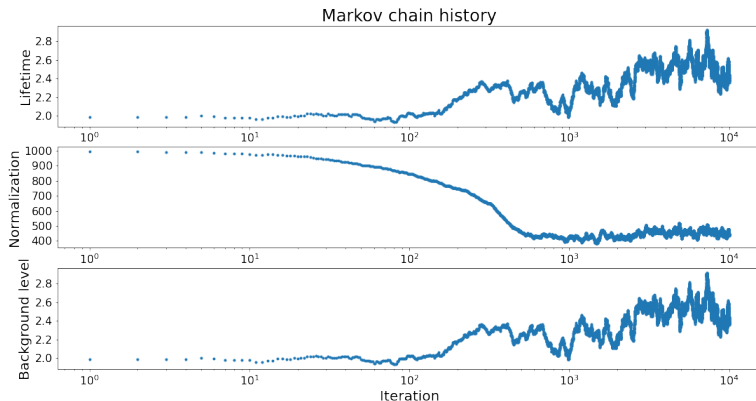
Markov Chain Monte Carlo does not work “out of the box”



It converges fast with the proper choice of parameter variation steps

Example

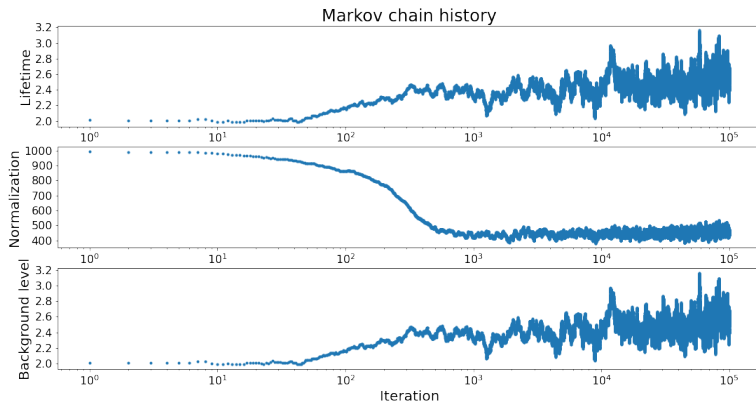
Markov Chain Monte Carlo does not work “out of the box”



Convergence can be very slow, if parameter steps too small...

Example

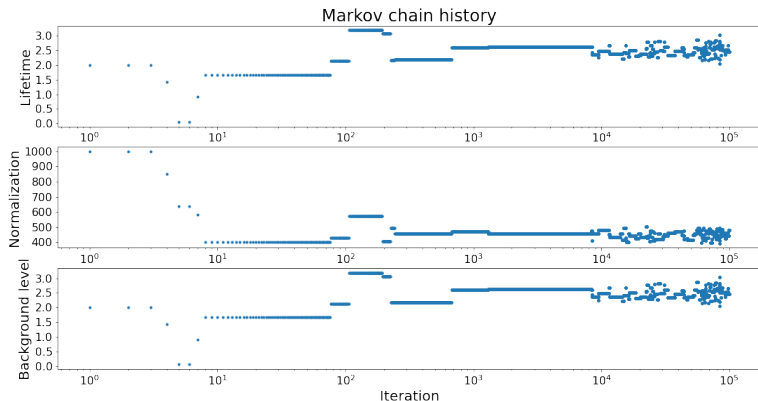
Markov Chain Monte Carlo does not work “out of the box”



Convergence can be very slow, if parameter steps too small...

Example

Markov Chain Monte Carlo does not work “out of the box”



Fluctuations significantly increased, if steps are too large...

Final remarks

Markov Chains are powerful tools to solve many problems that are difficult to approach “directly”, using other numerical techniques

However, it is crucial to make sure they converge, before using their output for the analysis. **Algorithm tuning may be required...**

Final remarks

Markov Chains are powerful tools to solve many problems that are difficult to approach “directly”, using other numerical techniques

However, it is crucial to make sure they converge, before using their output for the analysis. **Algorithm tuning may be required...**

Only the simplest approach was presented, many more advanced algorithms exist for more effective step generation

Probability $p(X^{(t+1)}|X^{(t)})$ does not need to be uniform!

Final remarks

Markov Chains are powerful tools to solve many problems that are difficult to approach “directly”, using other numerical techniques

However, it is crucial to make sure they converge, before using their output for the analysis. **Algorithm tuning may be required...**

Only the simplest approach was presented, many more advanced algorithms exist for more effective step generation

Probability $p(X^{(t+1)}|X^{(t)})$ does not need to be uniform!

Events generated with Markov Chain MC are not independent!

One should not use subsequent events together in the analysis
(eg. for background estimates)