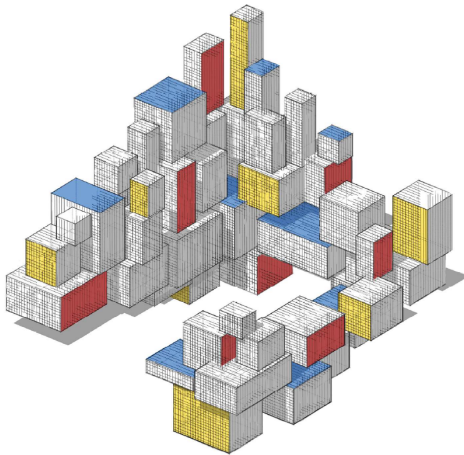


Markov Chain Monte Carlo



Purpose

In this lecture we discuss:

- Simulating Random Vectors
- Simulating Markov Chains
- Resampling Techniques
- Markov Chain Monte Carlo
 - Metropolis–Hastings Sampler
 - Gibbs Sampler

Simulating Random Vectors and Processes

Two general scenarios for simulating random vectors and processes:

1. X_1, \dots, X_n are **independent**, with pdfs $f_i, i = 1, \dots, n$. The random vector $\mathbf{X} = [X_1, \dots, X_n]^\top$ can be simply simulated by drawing each component $X_i \sim f_i$ individually — for example, via the inverse-transform method or acceptance–rejection.
2. For **dependent** components X_1, \dots, X_n , we can, as a consequence of the *product rule* of probability, represent the joint pdf $f(\mathbf{x})$ as

$$f(\mathbf{x}) = f(x_1, \dots, x_n) = f_1(x_1) f_2(x_2 | x_1) \cdots f_n(x_n | x_1, \dots, x_{n-1}).$$

One can generate \mathbf{X} by first generating X_1 , then, given $X_1 = x_1$, generate X_2 from $f_2(x_2 | x_1)$, and so on, until generating X_n from $f_n(x_n | x_1, \dots, x_{n-1})$.

Simulating Markov Chains

A **Markov chain** is a stochastic process $\{X_t, t = 0, 1, 2, \dots\}$ that satisfies the *Markov property*; meaning that for all t and s the conditional distribution of X_{t+s} given $X_u, u \leq t$, is the same as that of X_{t+s} given only X_t .

As a result, each conditional density $f_t(x_t | x_1, \dots, x_{t-1})$ can be written as a one-step **transition density** $q_t(x_t | x_{t-1})$; that is, the probability density to go from state x to state y in one step.

In many cases of interest the chain is **time-homogeneous**, meaning that the transition density q_t does not depend on t . Such Markov chains can be generated *sequentially*, as given in the following algorithm.

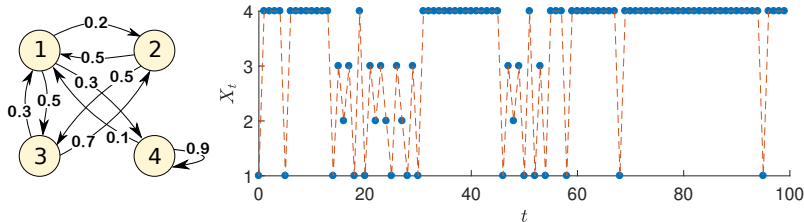
Algorithm 1: Simulate a Markov Chain

input: Number of steps N , initial pdf f_0 , transition density q .

- 1 Draw X_0 from the initial pdf f_0 .
 - 2 **for** $t = 1$ **to** N **do**
 - 3 Draw X_t from the distribution corresponding to the transition density $q(\cdot \mid X_{t-1})$
 - 4 **return** X_0, \dots, X_N
-

Example: Markov Chain Simulation

For time-homogeneous Markov chains with a discrete state space, we can visualize the one-step transitions by means of a **transition graph**, where arrows indicate possible transitions between states and the labels describe the corresponding probabilities. The following figure shows (on the left) the transition graph of the Markov chain $\{X_t, t = 0, 1, 2, \dots\}$ with state space $\{1, 2, 3, 4\}$, starting at state 1, and on the right a typical path of the Markov chain.



```
import numpy as np
import matplotlib.pyplot as plt

n = 101
P = np.array([[0, 0.2, 0.5, 0.3],
              [0.5, 0, 0.5, 0],
              [0.3, 0.7, 0, 0],
              [0.1, 0, 0, 0.9]])
x = np.array(np.ones(n, dtype=int))
x[0] = 0
for t in range(0,n-1):
    x[t+1] = np.min(np.where(np.cumsum(P[x[t],:]) >
                              np.random.rand()))
x = x + 1 #add 1 to all elements of the vector x
plt.plot(np.array(range(0,n)),x, 'o')
plt.plot(np.array(range(0,n)),x, '--')
plt.show()
```

Resampling

The idea behind **resampling** is very simple: an iid sample $\tau := \{x_1, \dots, x_n\}$ from some unknown cdf F represents our best knowledge of F if we make no further *a priori* assumptions about it. If it is not possible to simulate more samples from F , the best way to “repeat” the experiment is to resample each x_i with equal probability and repeat this N times.

Algorithm 2: Sampling from an Empirical Cdf.

input: Original iid sample x_1, \dots, x_n and sample size N .

output : Iid sample X_1^*, \dots, X_N^* from the empirical cdf.

```
1 for  $t = 1$  to  $N$  do
2   |   Draw  $U \sim \mathcal{U}(0, 1)$ 
3   |   Set  $I \leftarrow \lceil nU \rceil$       // discrete uniformly from  $1, \dots, n$ 
4   |   Set  $X_t^* \leftarrow x_I$ 
5 return  $X_1^*, \dots, X_N^*$ 
```

Example: Quotient of Uniforms

Let $U_1, \dots, U_n, V_1, \dots, V_n$ be iid $\mathcal{U}(0, 1)$ random variables and define $X_i = U_i/V_i, i = 1, \dots, n$.

Suppose we wish to investigate the distribution of the sample median \tilde{X} and sample mean \bar{X} of the (random) data $\mathcal{T} := \{X_1, \dots, X_n\}$.

Since we know the model for \mathcal{T} exactly, we can generate a large number, N say, of independent copies of it, and for each of these copies evaluate the sample medians $\tilde{X}_1, \dots, \tilde{X}_N$ and sample means $\bar{X}_1, \dots, \bar{X}_N$.

For $n = 100$ and $N = 1000$ the empirical cdfs might look like the left and right curves in Figure 1, respectively.

The distributions of the sample median and sample mean do not match at all. The sample median is quite concentrated around 1, whereas the distribution of the sample mean is much more spread out.

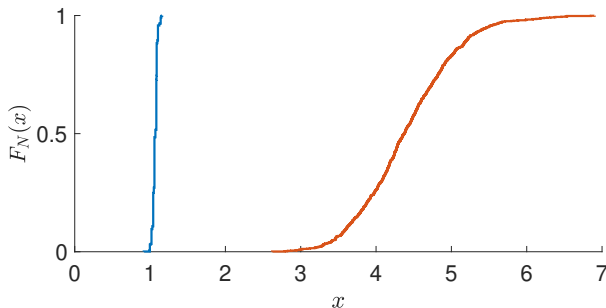


Figure: Empirical cdfs of the medians of the resampled data (left curve) and sample means (right curve) of the resampled data.

Instead of sampling completely new data, we could also *reuse* the original data by resampling them via Algorithm 2. This gives independent copies $\widetilde{X}_1^*, \dots, \widetilde{X}_N^*$ and $\overline{X}_1^*, \dots, \overline{X}_N^*$, for which we can again plot the empirical cdf.

The results will be similar to the previous case. In fact, in the figure above, the empirical cdfs of the *resampled* sample medians and sample means were plotted.

The corresponding Python code is given below. The essential point of this example is that resampling of data can greatly add to the understanding of the probabilistic properties of certain measurements on the data, *even if the underlying model is not known*.

```
import numpy as np
from numpy.random import rand, choice
import matplotlib.pyplot as plt
from statsmodels.distributions.empirical_distribution import ECDF

n = 100
N = 1000
x = rand(n)/rand(n) # data
med = np.zeros(N)
ave = np.zeros(N)
for i in range(0,N):
    s = choice(x, n, replace=True) # resampled data
    med[i] = np.median(s)
    ave[i] = np.mean(s)

med_cdf = ECDF(med)
ave_cdf = ECDF(ave)
plt.plot(med_cdf.x, med_cdf.y)
plt.plot(ave_cdf.x, ave_cdf.y)
plt.show()
```

Markov Chain Monte Carlo

Markov chain Monte Carlo (MCMC) is a Monte Carlo sampling technique for (approximately) generating samples from an arbitrary distribution — often referred to as the **target** distribution.

The basic idea is to run a Markov chain long enough such that its **limiting distribution** is close to the target distribution.

The initial random variables in the Markov chain may have a distribution that is significantly different from the target (limiting) distribution. The random variables that are generated during this **burn-in period** are often discarded.

The remaining random variables form an *approximate* and *dependent* sample from the target distribution.

We discuss two popular MCMC samplers: the Metropolis–Hastings sampler and the Gibbs sampler.

Metropolis–Hastings Sampler

We wish to sample from a target pdf $f(\mathbf{x})$, where \mathbf{x} takes values in some d -dimensional set. The aim is to construct a Markov chain $\{\mathbf{X}_t, t = 0, 1, \dots\}$ in such a way that its limiting pdf is f .

Suppose the Markov chain is in state \mathbf{x} at time t . A transition of the Markov chain from state \mathbf{x} is carried out in two phases. First a **proposal** state \mathbf{y} is drawn from a transition density $q(\cdot | \mathbf{x})$. This state is accepted as the new state, with **acceptance probability**

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{f(\mathbf{y}) q(\mathbf{x} | \mathbf{y})}{f(\mathbf{x}) q(\mathbf{y} | \mathbf{x})}, 1 \right\}, \quad (1)$$

or rejected otherwise. In the latter case the chain remains in state \mathbf{x} . The algorithm just described can be summarized as follows.

Algorithm 3: Metropolis–Hastings Sampler

input: Initial state X_0 , sample size N , target pdf $f(\mathbf{x})$, proposal function $q(\mathbf{x}, \mathbf{y})$.

output: X_1, \dots, X_N (dependent), approximately distributed according to $f(\mathbf{x})$.

```
1 for  $t = 0$  to  $N - 1$  do
2   Draw  $Y \sim q(\mathbf{y} \mid X_t)$  // draw a proposal
3    $\alpha \leftarrow \alpha(X_t, Y)$  // acceptance probability as in (1)
4   Draw  $U \sim \mathcal{U}(0, 1)$ 
5   if  $U \leq \alpha$  then  $X_{t+1} \leftarrow Y$ 
6   else  $X_{t+1} \leftarrow X_t$ 
7 return  $X_1, \dots, X_N$ 
```

Note that one only needs to know the target pdf $f(\mathbf{x})$ *up to a constant*; that is $f(\mathbf{x}) = c \bar{f}(\mathbf{x})$ for some known function $\bar{f}(\mathbf{x})$ but unknown constant c .

Detailed Balance Equations

Because a transition of the Metropolis–Hastings Markov chain consists of two steps, the one-step transition probability to go from \mathbf{x} to \mathbf{y} is not $q(\mathbf{y} | \mathbf{x})$ but

$$\tilde{q}(\mathbf{y} | \mathbf{x}) = \begin{cases} q(\mathbf{y} | \mathbf{x}) \alpha(\mathbf{x}, \mathbf{y}), & \text{if } \mathbf{y} \neq \mathbf{x}, \\ 1 - \sum_{\mathbf{z} \neq \mathbf{x}} q(\mathbf{z} | \mathbf{x}) \alpha(\mathbf{x}, \mathbf{z}), & \text{if } \mathbf{y} = \mathbf{x}. \end{cases}$$

This transition density satisfies the **detailed balance equations**:

$$f(\mathbf{x}) \tilde{q}(\mathbf{y} | \mathbf{x}) = f(\mathbf{y}) \tilde{q}(\mathbf{x} | \mathbf{y}) \quad \text{for all } \mathbf{x}, \mathbf{y}.$$

This is the reason why the limiting distribution of the Metropolis–Hastings Markov chain is equal to the target distribution.

Choice of the Proposal

Ideally, we would like $q(\mathbf{y} | \mathbf{x})$ to be “close” to the target $f(\mathbf{y})$, irrespective of \mathbf{x} . Two common approaches are:

1. **Independence Sampler:** Choose $q(\mathbf{y} | \mathbf{x}) = g(\mathbf{y})$ independent of \mathbf{x} for some pdf $g(\mathbf{y})$. The acceptance probability is thus

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{f(\mathbf{y}) g(\mathbf{x})}{f(\mathbf{x}) g(\mathbf{y})}, 1 \right\}.$$

2. **Random Walk Sampler:** Choose a symmetric proposal density $q(\mathbf{y} | \mathbf{x}) = q(\mathbf{x} | \mathbf{y})$. Then the acceptance probability has the simple form

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{f(\mathbf{y})}{f(\mathbf{x})}, 1 \right\}. \quad (2)$$

A typical example is when, for a given current state \mathbf{x} , the proposal state \mathbf{Y} is of the form $\mathbf{Y} = \mathbf{x} + \mathbf{Z}$, where \mathbf{Z} is generated from some spherically symmetric distribution, such as $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

Example: Random Walk Sampler

For an unspecified constant c , consider the two-dimensional pdf

$$f(x_1, x_2) = c e^{-\frac{1}{4}\sqrt{x_1^2 + x_2^2}} \left(\sin \left(2\sqrt{x_1^2 + x_2^2} \right) + 1 \right), \quad [x_1, x_2] \in [-2\pi, 2\pi]^2.$$

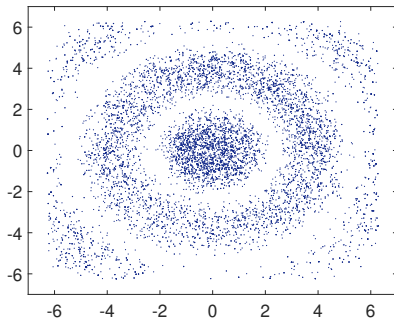
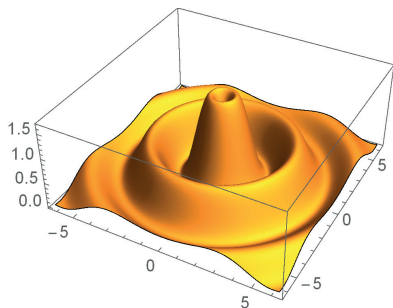


Figure: Left panel: the two-dimensional target pdf. Right panel: points from the random walk sampler. At each step, given a current state \mathbf{x} , a proposal \mathbf{Y} is drawn from the $\mathcal{N}(\mathbf{x}, \mathbf{I})$ distribution.

```
import numpy as np
import matplotlib.pyplot as plt
from numpy import pi, exp, sqrt, sin
from numpy.random import rand, randn
N = 10000
a = lambda x: -2*pi < x
b = lambda x: x < 2*pi
f = lambda x1, x2: exp(-sqrt(x1**2+x2**2)/4)*(
    sin(2*sqrt(x1**2+x2**2))+1)*a(x1)*b(x1)*a(x2)*b(x2)

xx = np.zeros((N,2))
x = np.zeros((1,2))
for i in range(1,N):
    y = x + randn(1,2)
    alpha = np.amin((f(y[0][0],y[0][1])/f(x[0][0],x[0][1]),1))
    r = rand() < alpha
    x = r*y + (1-r)*x
    xx[i,:] = x

plt.scatter(xx[:,0], xx[:,1], alpha=0.4,s=2)
plt.axis('equal')
```

Gibbs Sampler

The **Gibbs sampler** is useful for generating n -dimensional random vectors. The key idea of the Gibbs sampler is to update the components of the random vector one at a time, by sampling them from *conditional* pdfs.

Specifically, suppose that we wish to sample a random vector $\mathbf{X} = [X_1, \dots, X_n]^\top$ according to a target pdf $f(\mathbf{x})$.

Using Bayesian notation, let $f(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ represent the conditional pdf of the i -th component, X_i , given the other components $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$.

Systematic Gibbs Sampler

Algorithm 4: Systematic Gibbs Sampler

input: Initial point X_0 , sample size N , and target pdf f .

output: X_1, \dots, X_N approximately distributed according to f .

```
1 for  $t = 0$  to  $N - 1$  do
2   Draw  $Y_1$  from the conditional pdf  $f(y_1 \mid X_{t,2}, \dots, X_{t,n})$ .
3   for  $i = 2$  to  $n$  do
4     Draw  $Y_i$  from the conditional pdf
        $f(y_i \mid Y_1, \dots, Y_{i-1}, X_{t,i+1}, \dots, X_{t,n})$ .
5    $X_{t+1} \leftarrow Y$ 
6 return  $X_1, \dots, X_N$ 
```

Note that here the components are updated in a fixed order $1 \rightarrow 2 \rightarrow \dots \rightarrow n$. For this reason Algorithm 4 is also called the **systematic Gibbs sampler**.

Variants of the Gibbs Sampler

There exist many variants of the Gibbs sampler, depending on **cycle** steps required to update \mathbf{X}_t to \mathbf{X}_{t+1} .

- **Random-order Gibbs sampler**: update the components in the order of a random permutation of $\{1, \dots, n\}$.
- Update the components in blocks (i.e., several at the same time).
- Update only a random selection of components; e.g., only a single component (**random Gibbs sampler**).
- **Reversible Gibbs sampler**: update according to $1 \rightarrow 2 \rightarrow \dots \rightarrow n-1 \rightarrow n \rightarrow n-1 \rightarrow \dots \rightarrow 2 \rightarrow 1$.

In all cases, except for the systematic Gibbs sampler, the resulting Markov chain $\{\mathbf{X}_t, t = 1, 2, \dots\}$ is *reversible* and hence its limiting distribution is precisely $f(\mathbf{x})$.

Instead, the systematic Gibbs sampler satisfies a related balance condition (Hammersley–Clifford), which guarantees that the limiting pdf is indeed $f(\mathbf{x})$.

Example: Gibbs Sampler for the Bayesian Normal Model

Gibbs samplers are often applied in Bayesian statistics, to sample from the posterior pdf. Consider for instance the Bayesian normal model

$$f(\mu, \sigma^2) = 1/\sigma^2$$
$$(\mathbf{x} \mid \mu, \sigma^2) \sim \mathcal{N}(\mu \mathbf{1}, \sigma^2 \mathbf{I}).$$

The prior for (μ, σ^2) is here [improper](#). In some sense this prior conveys the least amount of information about μ and σ^2 .

Following the usual procedure, we find the posterior pdf:

$$f(\mu, \sigma^2 \mid \mathbf{x}) \propto (\sigma^2)^{-n/2-1} \exp \left\{ -\frac{1}{2} \frac{\sum_i (x_i - \mu)^2}{\sigma^2} \right\}. \quad (3)$$

Note that μ and σ^2 here are the “variables” and \mathbf{x} is a fixed data vector.

Simulate from the Posterior

To simulate samples μ and σ^2 from (3) using the Gibbs sampler, we need the distributions of both $(\mu \mid \sigma^2, \mathbf{x})$ and $(\sigma^2 \mid \mu, \mathbf{x})$.

To find $f(\mu \mid \sigma^2, \mathbf{x})$, view the right-hand side of (3) as a function of μ only, regarding σ^2 as a constant. This gives

$$\begin{aligned} f(\mu \mid \sigma^2, \mathbf{x}) &\propto \exp \left\{ -\frac{n\mu^2 - 2\mu \sum_i x_i}{2\sigma^2} \right\} = \exp \left\{ -\frac{\mu^2 - 2\mu\bar{x}}{2(\sigma^2/n)} \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \frac{(\mu - \bar{x})^2}{\sigma^2/n} \right\}. \end{aligned} \quad (4)$$

This shows that $(\mu \mid \sigma^2, \mathbf{x})$ has a normal distribution with mean \bar{x} and variance σ^2/n .

Similarly, to find $f(\sigma^2 | \mu, \mathbf{x})$, view the right-hand side of (3) as a function of σ^2 , regarding μ as a constant. This gives

$$f(\sigma^2 | \mu, \mathbf{x}) \propto (\sigma^2)^{-n/2-1} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 / \sigma^2 \right\}, \quad (5)$$

showing that $(\sigma^2 | \mu, \mathbf{x})$ has an inverse-gamma distribution with parameters $n/2$ and $\sum_{i=1}^n (x_i - \mu)^2 / 2$. The Gibbs sampler thus

involves the repeated simulation of

- $(\mu | \sigma^2, \mathbf{x}) \sim \mathcal{N}(\bar{x}, \sigma^2/n)$ and
- $(\sigma^2 | \mu, \mathbf{x}) \sim \text{InvGamma}(n/2, \sum_{i=1}^n (x_i - \mu)^2 / 2)$.

Simulating $X \sim \text{InvGamma}(\alpha, \lambda)$ is achieved by first generating $Z \sim \text{Gamma}(\alpha, \lambda)$ and then returning $X = 1/Z$.

```
import numpy as np
import matplotlib.pyplot as plt

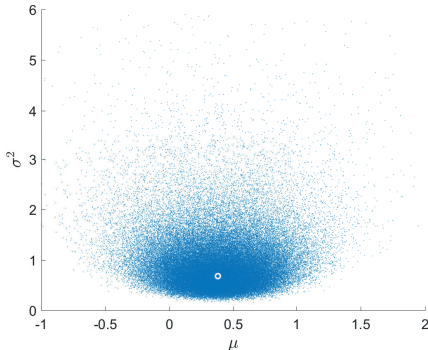
x = np.array([[ -0.9472, 0.5401, -0.2166, 1.1890, 1.3170,
                -0.4056, -0.4449, 1.3284, 0.8338, 0.6044]])

n=x.size
sample_mean = np.mean(x)
sample_var = np.var(x)
sig2 = np.var(x)
mu=sample_mean

N=10**5
gibbs_sample = np.array(np.zeros((N, 2)))
for k in range(N):
    mu=sample_mean + np.sqrt(sig2/n)*np.random.randn()
    V=np.sum((x-mu)**2)/2
    sig2 = 1/np.random.gamma(n/2, 1/V)
    gibbs_sample[k,:]= np.array([mu, sig2])
plt.scatter(gibbs_sample[:,0], gibbs_sample[:,1],alpha =0.1,s =1)
plt.plot(np.mean(x), np.var(x), 'wo')
```

Output from the Gibbs Sampler

The figure shows the posterior point cloud; i.e., the (μ, σ^2) points generated by the Gibbs sampler. The white circle is the point (\bar{x}, s^2) , where $\bar{x} = 0.3798$ is the sample mean and $s^2 = 0.6810$ the sample variance.



Posterior for μ

By projecting the (μ, σ^2) points onto the μ -axis — that is, by ignoring the σ^2 values — one obtains (approximate) samples from the posterior pdf of μ ; that is, $f(\mu | \mathbf{x})$.

The figure shows a kernel density estimate of this pdf.

The corresponding 0.05 and 0.95 sample quantiles were found to be -0.2054 and 0.9662 , respectively, giving the 95% credible interval $(-0.2054, 0.9662)$ for μ , which contains the true expectation 0.

