

Time Series Analysis

6. Markov Chain Monte Carlo

Andrew Lesniewski

Baruch College
New York

Fall 2019

Outline

- 1 Bayesian inference versus classical inference
- 2 Markov chains
- 3 Metropolis-Hastings algorithm

Bayesian viewpoint on inference

- Let x denote a single observation of a random variable (data point); this may be a number, a vector, or a more complex object.
- By θ , we denote the parameters of the probability distribution x . As usual, we use the notation

$$x \sim p(x|\theta), \quad (1)$$

to indicate that x is distributed according to the law $p(x|\theta)$.

- The objective of statistical inference is to predict the probability distribution of a new data point x , given a set x_1, \dots, x_N of N observed data points.
- In the previous lectures we took the classical (“frequentist”) point of view: in order to predict the PDF $p(x|\theta)$, we *plug* into $p(x|\theta)$ some optimized value of θ (for example the MLE estimate $\hat{\theta}$ of θ).

Bayesian viewpoint on inference

- This classical paradigm to estimating θ is explicitly based on the assumption that all probabilities involved have an *objective* character, and the MLE estimator $\hat{\theta}$ gives the best fit based on the observations.
- The objectivity assumption is important in some scientific inquiries, where the reproducibility of results is a key criterion of their validity.
- Such an objectivity requirement is, however, impractical (or outright impossible) in finance, where market outcomes are unique and irreproducible.
- Probability distributions are based on the participants' subjective views, or reflect current market sentiment as implied from option prices.
- A more general and flexible inference paradigm is *Bayesian inference*, which naturally integrates the concept of *subjective probability*.

Bayesian viewpoint on inference

- Bayesian inference is based on Bayes' definition of conditional probability,

$$P(E|H) = \frac{P(E \cap H)}{P(H)}. \quad (2)$$

- One can think about H as a “hypothesis” and about E as “evidence”.
- Interchanging the roles of E and H ,

$$P(H|E) = \frac{P(E \cap H)}{P(E)},$$

and equating the two expressions for $P(E \cap H)$, we obtain the following *Bayes' rule*:

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}. \quad (3)$$

Bayesian viewpoint on inference

- Relation (3) is phrased as

$$\text{posterior probability} = \frac{\text{likelihood} \times \text{prior probability}}{\text{evidence}}. \quad (4)$$

Since the denominator is a constant independent of H , this is often stated as a proportionality

$$\text{posterior probability} \propto \text{likelihood} \times \text{prior probability}. \quad (5)$$

- The prior probability $P(H)$ is a (subjective) probability that the hypothesis holds, while the likelihood $P(E|H)$ is the probability of E occurring given H .
- The *prior distribution* is the (subjective) probability distribution of the parameters, representing one's beliefs about them, before any evidence is observed. It is denoted by $p(\theta)$.

Bayesian viewpoint on inference

- The likelihood function is the distribution of the observed data conditioned on its parameters, i.e. $p(x|\theta)$. We use the notation $\mathcal{L}(\theta|x) = p(x|\theta)$ to emphasize that it is a function of the parameters.
- The *marginal likelihood* (a.k.a. the *evidence*) is the distribution of the observed data, i.e. the marginal PDF over the parameters θ ,

$$p(x) = \int p(x|\theta)p(\theta)d\theta. \quad (6)$$

- The *posterior distribution* is the distribution of the parameters after taking into account the evidence. From Bayes' rule (3),

$$\begin{aligned} p(\theta|x) &= \frac{p(x|\theta)p(\theta)}{p(x)} \\ &\propto p(x|\theta)p(\theta). \end{aligned} \quad (7)$$

- This identity is the essence of Bayesian inference.

Bayesian viewpoint on inference

- The *posterior predictive distribution* is the PDF of a new data point x' given the posterior:

$$p(x'|x) = \int p(x'|\theta)p(\theta|x)d\theta. \quad (8)$$

- According to the strict Bayesian paradigm one has to use the posterior predictive distribution in order to predict the distribution of a new, unobserved data point.
- In other words, an entire probability distribution of the parameters θ should be found rather than a single value of θ .
- This is frequently impractical, and short cuts have to be taken.
- Prediction in classical inference involves finding a single optimum point estimate of the parameters, namely the MLE estimator:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \mathcal{L}(\theta|x). \quad (9)$$

Bayesian viewpoint on inference

- A half way approach is the *maximum a posteriori estimation* (MAP), in which the posterior distribution is maximized. Namely,

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} p(\theta|x) \\ &= \arg \max_{\theta} \frac{p(x|\theta)p(\theta)}{p(x)} \\ &= \arg \max_{\theta} p(x|\theta)p(\theta).\end{aligned}\tag{10}$$

- As an example, consider a linear regression model

$$y = \alpha + \beta x,\tag{11}$$

given the observations $(x_1, y_1), \dots, (x_N, y_N)$. Assuming normality, the likelihood function of this data set is

$$p(x, y|\theta) = (2\pi\sigma^2)^{-N/2} \prod_{j=1}^N \exp\left(-\frac{(y_j - \alpha - \beta x_j)^2}{2\sigma^2}\right).$$

Bayesian viewpoint on inference

- Maximizing this likelihood function leads to the standard regression estimators for the parameters α and β :

$$\begin{aligned}\hat{\beta} &= \frac{\sum_j (y_j - \hat{y})(x_j - \hat{x})}{\sum_j (x_j - \hat{x})^2}, \\ \hat{\alpha} &= \hat{y} - \hat{\beta} \hat{x},\end{aligned}\tag{12}$$

where $\hat{y} = \frac{1}{N} \sum_j y_j$, and $\hat{x} = \frac{1}{N} \sum_j x_j$.

- On the other hand, we might want to express the view that the parameter β should be close to zero. This can be done by selecting a Gaussian prior

$$p(\beta) \propto \exp\left(-\frac{\lambda\beta^2}{2}\right),$$

where $\lambda > 0$ is a suitably chosen *hyperparameter*.

Bayesian viewpoint on inference

- MAP estimation leads then to the problem of minimizing the following posterior function:

$$-\log p(\theta|x, y) = \frac{N}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{j=1}^N (y_j - \alpha - \beta x_j)^2 + \frac{\lambda \beta^2}{2}.$$

- The optimal values of the coefficients are

$$\hat{\beta} = \frac{\sum_j (y_j - \hat{y})(x_j - \hat{x})}{\sum_j (x_j - \hat{x})^2 + \lambda}, \quad (13)$$
$$\hat{\alpha} = \hat{y} - \hat{\beta} \hat{x}.$$

- This results in the *ridge regression*, in which the parameter β is shrunk towards zero by λ .

Bayesian viewpoint on inference

- Likewise, selecting the prior

$$p(\beta) \propto e^{-\lambda|\beta|},$$

with $\lambda > 0$, leads to the *lasso regression*, where the objective is to minimize the posterior function:

$$-\log p(\theta|x, y) = \frac{N}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{j=1}^N (y_j - \alpha - \beta x_j)^2 + \lambda|\beta|.$$

- The only difference between the MLE and MAP estimates is the presence of the prior probability: the MAP estimate allows us to inject our prior beliefs regarding the values of θ into the calculation .
- The classical approach has the disadvantage that it does not account for any uncertainty in the values of the parameters, and thus tends to underestimate the variance of the predictive distribution.

Conjugate priors

- Full Bayesian inference involves calculating the constant factor (6) in the denominator (the evidence).
- This is, in general, a difficult task (it involves multi-dimensional integration), and it makes full Bayesian inference a hard problem.
- One way to address this issue is to, for a given likelihood function, express our prior beliefs in a way that allow us to carry out the integration in (6).
- This leads to the notion of a *conjugate prior*.
- If the posterior distribution $p(\theta|x)$ is in the same family of distributions as the prior PDF $p(\theta)$, the prior and posterior are called conjugate distributions.
- In such a situation, the prior is called a conjugate prior for the likelihood function.

Conjugate priors

- As an example, assume that the likelihood function is normal with known standard deviation σ , i.e. the only model parameter is the mean μ .
- Taking the prior to be normal with (hyper)parameters μ_0, σ_0 , we find that the posterior is normal with mean

$$\mu' = \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{j=1}^N x_j \right) / \left(\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \right), \quad (14)$$

and standard deviation

$$\sigma' = 1 / \sqrt{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}. \quad (15)$$

Conjugate priors

- Conjugate priors are known explicitly only in a limited number of cases.
- Even then, the calculations involved may be quite formidable.
- A practical approach to Bayesian computations is through Monte Carlo methods, and in particular through Markov chain Monte Carlo methods.

Markov chains

- Consider a discrete time process θ_t , $t = 0, 1, \dots$, with values in a *state space* A , $\theta_t \in A$.
- A can be continuous (such as \mathbb{R} or \mathbb{R}^n) or discrete (such as $\{a_1, \dots, a_n\}$).
- The process θ_t is called a *Markov chain*, if

$$P(\theta_{t+1} | \theta_{1:t}) = P(\theta_{t+1} | \theta_t), \quad (16)$$

or, in other words, if θ_{t+1} is independent of θ_s with $s < t$.

- We will be considering Markov chains with *stationary transition probabilities*, i.e. such that $P(\theta_{t+1} | \theta_t)$ is independent of t .
- This is, in particular, true if the process is strictly stationary.

Markov chains

- If the state space is finite A , the transition probabilities can be encoded in the matrix $P \in \text{Mat}_n(\mathbb{R})$, defined by

$$p_{ij} = P(\theta_{t+1} = a_j | \theta_t = a_i). \quad (17)$$

- Notice that, for all $i = 1, \dots, n$,

$$\sum_{j=1}^n p_{ij} = 1. \quad (18)$$

- Matrices with non-negative elements satisfying the above identity are called *stochastic matrices*.

Markov chains

- The initial probability distribution π_0 is a vector defined as

$$\pi_{0,j} = \mathbf{P}(\theta_0 = a_j). \quad (19)$$

- The *marginal probability distribution* π_t at time $t = 1, 2, \dots$, has components

$$\pi_{t,i} = \sum_{j=1}^n p_{ji} \pi_{t-1,j}.$$

- In matrix notation this reads

$$\pi_t = \pi_{t-1} P, \quad (20)$$

or

$$\pi_t = \pi_0 P^t. \quad (21)$$

Equilibrium and detailed balance

- The marginal probability distribution π is called an *equilibrium distribution*, if $\pi_t = \pi_{t-1}$, i.e.

$$\pi = \pi P. \quad (22)$$

- The notion of an equilibrium distribution is closely related to the *detailed balance condition*, namely

$$\pi_i p_{ij} = \pi_j p_{ji}. \quad (23)$$

- A Markov chain for which there is a probability distribution π satisfying the detailed balance condition is called *reversible*.
- This terminology is justified by the observation that (23) states that the absolute probabilities of moving from state i to j and from j to i are identical.
- In other words, the Markov chain remains invariant under reversing the time direction.

Equilibrium and detailed balance

- We claim that *if π satisfies the detailed balance condition, then it is an equilibrium distribution.*
- The proof is a straightforward calculation:

$$\begin{aligned}(\pi P)_j &= \sum_{i=1}^n \pi_i p_{ij} \\ &= \sum_{i=1}^n \pi_j p_{ji} \\ &= \pi_j \sum_{i=1}^n p_{ji} \\ &= \pi_j,\end{aligned}$$

i.e. $\pi P = \pi$. QED

Equilibrium and detailed balance

- It is, however, in general not true that an equilibrium distribution satisfies the detailed balance condition.
- The concepts above can be generalized to the case of a continuous state space by replacing the sums with suitable integrals.

Convergence to equilibrium

- An important practical question is whether a given Markov chain has an equilibrium distribution and whether it is unique.
- We assume the the Markov chain is irreducible, i.e. it is possible to reach from one state to any other state with non-zero probability.
- Consider the limit

$$P_\infty = \lim_{t \rightarrow \infty} P^t. \quad (24)$$

Then $P_\infty P = P_\infty$, and all rows of P_∞ are equal. We will denote this common value by π .

- In other words, π represents an equilibrium distribution of the Markov chain. When is it unique?

Convergence to equilibrium

- A state θ_j is called *ergodic*, if it is
 - (i) recurrent, i.e. the system returns to this state with probability 1, and the expected time of the return is finite;
 - (ii) aperiodic, i.e. the recurrence is non-periodic.

If all states in an irreducible Markov chain are ergodic, then the chain is called ergodic.

- **Theorem 1.** *For an ergodic Markov chain, there is a unique probability vector π such that $\pi P = \pi$.*
- **Theorem 2.** *For an ergodic Markov chain,*

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T P^t = P_{\infty}. \quad (25)$$

In other words, in an ergodic Markov chain, the time average of a state is equal to its average under the equilibrium probability distribution.

- The rate of convergence to the equilibrium distribution depends on the details of the transition matrix.

Markov chain Monte Carlo

- *Markov chain Monte Carlo* (MCMC) methods are methodologies for generating samples from a probability distribution $p(\theta)$ using the techniques of Markov chains.
- The underlying idea is to construct a Markov chain that has the desired distribution $p(\theta)$ as its equilibrium distribution. The states of the chain are then used as samples from the distribution $p(\theta)$.
- In applications, the distribution $p(\theta)$ may not be given by an explicit formula, or may be otherwise hard to sample from.
- MCMC methods are used for calculating numerical approximations of multi-dimensional integrals, which arise in Bayesian statistics.
- This is accomplished by generating random samples from the prior or posterior distributions.
- Such applications may require integrations over hundreds or thousands of unknown parameters.

Classic Monte Carlo method

- MCMC methods should be contrasted with the classic Monte Carlo methods used in many financial applications (such as valuation models).
- Classic Monte Carlo utilizes non-Markov chain simulated random variables.
- The basic input into classic Monte Carlo methods is an efficient algorithm for generating samples from the uniform distribution $U(0, 1)$ on the interval $(0, 1)$.
- These samples can then be transformed into samples from the standard normal distribution, as well as a number of other parametric distributions.
- Typically, the simulated random variables are statistically independent; any desired dependence is artificially introduced using e.g. the Cholesky decomposition of the covariance matrix.
- Simulation from any given distribution $p(\theta)$ may pose challenges as no simple transformation onto the uniform distribution is possible.

Acceptance-rejection method

- A general technique for simulation of independent random variables from a given PDF $p(\theta)$ is the *acceptance-rejection method*.
- The PDF $p(\theta)$ is usually given in the form:

$$p(\theta) = \frac{1}{Z} f(\theta). \quad (26)$$

- Here $f(\theta)$ is an “unnormalized” density function, and Z is the normalization constant,

$$Z = \int f(\theta) d\theta, \quad (27)$$

whose value may not be explicitly known.

Acceptance-rejection method

- Let $q(\theta)$ be another PDF, called the *proposal distribution*, with the properties that:
 - there is an effective method for simulating from $q(\theta)$,
 - there is a constant M (explicitly known) such that for all θ ,

$$f(\theta) \leq Mq(\theta). \quad (28)$$

- The acceptance-rejection method is a two step process. In order to generate a random number from $p(\theta)$,
 - generate a candidate $\eta \sim q(\theta)$, and $u \sim U(0, 1)$ (the uniform distribution on $(0, 1)$),
 - if

$$u \leq \frac{f(\eta)}{Mq(\eta)}, \quad (29)$$

then accept η , else reject it and go back to (i).

Acceptance-rejection method

- Drawing a graph, we see that the algorithm indeed generates samples from $p(\theta)$.
- It is helpful to choose the constant M as close to the optimal as possible; large M slows down the algorithm.
- Even with the optimal choice of M , the algorithm may lead to a large number of rejections and poor performance.
- This is especially true in the case the proposal PDF significantly differs from the target PDF.

Markov chain Monte Carlo

- There are various approaches to MCMC. Two most popular approaches are:
 - (i) Metropolis-Hastings algorithm,
 - (ii) Gibbs sampler.
- These algorithms are clever extensions of the A-R algorithm by allowing the proposal distribution to depend on the state of a Markov chain.
- The Markov chain is selected so that the distribution $p(\theta)$ is its equilibrium distribution.
- The objective of the algorithm is to generate samples from a (complicated) target distribution $p(\theta)$. We denote its unnormalized density by $f(\theta)$, i.e. $p(\theta) \sim f(\theta)$.

Formulation of the algorithm

- Let $q(\eta|\theta)$ be a *proposal* PDF, with $\int q(\eta|\theta)d\eta = 1$.
 - (i) We interpret the proposal distribution by saying that if the chain is in the state θ , then it generates random variables (“candidates”) η .
 - (ii) The proposal PDF is selected so that it is easy to sample from it.
- Consider first the lucky situation, when $q(\theta|\eta)$ satisfies the detailed balance condition,

$$p(\theta)q(\eta|\theta) = p(\eta)q(\theta|\eta), \quad (30)$$

for all θ .

- Then

$$\begin{aligned} \sum_{\theta} p(\theta)q(\eta|\theta) &= p(\eta) \sum_{\theta} q(\theta|\eta) \\ &= p(\eta), \end{aligned}$$

and so a candidate η from $q(\eta|\theta)$ would also be a sample from $p(\eta)$, and it would always be accepted.

Formulation of the algorithm

- In general though, the detailed balance condition will be violated.
- For example, the process may be more likely to move from θ to η , i.e.

$$p(\theta)q(\eta|\theta) > p(\eta)q(\theta|\eta). \quad (31)$$

- This can be “corrected” by introducing the *probability of a move* $\alpha(\theta, \eta)$ from θ to η , so that the total transition probability is equal to

$$\tilde{p}(\eta|\theta) = q(\eta|\theta)\alpha(\theta, \eta). \quad (32)$$

Formulation of the algorithm

- In contrast to the A-R algorithm, if the move is not made, the system returns to the original state, and returns θ as the sample from $p(\theta)$.
- Let us now determine $\alpha(\theta, \eta)$. We will do it by requiring that $\tilde{p}(\eta|\theta)$ satisfies the detailed balance condition:

$$p(\theta)\tilde{p}(\eta|\theta) = p(\eta)\tilde{p}(\theta|\eta),$$

or

$$\begin{aligned} p(\theta)q(\eta|\theta)\alpha(\theta, \eta) &= p(\eta)q(\theta|\eta)\alpha(\eta, \theta) \\ &= p(\eta)q(\theta|\eta). \end{aligned} \tag{33}$$

- The last equality holds because, under the scenario (31), it is natural to require that $\alpha(\eta, \theta) = 1$.

Formulation of the algorithm

- This leads to the following value for $\alpha(\theta, \eta)$:

$$\alpha(\theta, \eta) = \frac{p(\eta)q(\theta|\eta)}{p(\theta)q(\eta|\theta)}.$$

- Notice that, importantly, the formula above does not depend on the normalizing factor in $p(\theta)$ and we can replace $p(\theta)$ by $f(\theta)$.
- The same analysis can be done if the sense of the inequality in (31) is reversed, leading to the definition of $\alpha(\eta, \theta)$.
- We summarize these calculations in the following manner. Let us form the *Hastings ratio*:

$$r(\theta, \eta) = \frac{f(\eta)q(\theta|\eta)}{f(\theta)q(\eta|\theta)}, \quad (34)$$

and define

$$\alpha(\theta, \eta) = \min(r(\theta, \eta), 1). \quad (35)$$

Formulation of the algorithm

- We are ready now ready to formulate the *Metropolis-Hastings (M-H) algorithm*.
- Choose an initial value θ_0 , and for $j = 1, 2, \dots$, and proceeds as follows:
 - (i) Generate a candidate $\eta \sim q(\eta|\theta_j)$ and $u \sim U(0, 1)$.
 - (ii) Accept η with probability $\alpha(\theta_j, \eta)$, i.e. accept if

$$u \leq \alpha(\theta_j, \eta), \quad (36)$$

and reject otherwise.

- (iii) If accepted, set $\theta_{j+1} = \eta$, else, $\theta_{j+1} = \theta_j$.
- The simulated values are $\theta_1, \theta_2, \dots$

Metropolis algorithm

- An important special case is when the proposal distribution is symmetric, $q(\theta|\eta) = q(\eta|\theta)$, in which case

$$r(\theta, \eta) = \frac{f(\eta)}{f(\theta)}. \quad (37)$$

- This case, historically preceding the general M-H algorithm is referred to as the *Metropolis algorithm*.

Burn-in period

- The simulated values can be considered as drawn from the PDF $p(x)$ only after an initial transient period.
- The length of the burn-in period depends on the rate at which the probability distribution of the Markov chain converges to its stationary limit.
- The random variables generated by the M-H algorithm are not necessarily independent of each other.
- For many purposes, such as evaluation of integrals, this is not a major concern.
- Other methodologies, such as *Hamiltonian MCMC* address this issue.

Choice of the proposal distribution

- A suitable choice of the proposal PDF is a key part of the implementation of the M-H algorithm.
- Typically, the proposal distribution is specified as a parameterized family of distributions, with tunable parameters.
- A possible choice of $q(\eta|\theta)$ is a suitable multivariate distribution $\varphi(\eta - \theta)$. In this case, the candidate is given by $\eta = \theta + \zeta$, where $\zeta \sim \varphi$.
- This variant of the algorithm is called the *random walk Metropolis*, because the candidate is chosen as the current state plus noise.
- If φ is symmetric around the origin, $\varphi(-\zeta) = \varphi(\zeta)$, then $q(\eta|\theta)$ is symmetric, and so the M-H algorithm reduces to the Metropolis algorithm.

Choice of the proposal distribution

- Popular choices of φ are:
 - (i) multivariate normal distribution,
 - (ii) multivariate t -distribution,
 - (iii) multivariate uniform distribution.
- Another common choice is $q(\eta|\theta) = \varphi(\eta)$, i.e. the candidates are drawn independently from a fixed distribution, independently of the current state.
- The resulting proposal distribution is non-symmetric and the full M-H algorithm has to be used.
- Again, the distribution φ can be specified as e.g. a (multivariate) normal distribution, t -distribution, gamma distribution, or uniform distribution.
- This version of the algorithm is known as the *independent M-H* algorithm.

Example: perturbed Gaussian distribution

- To illustrate the method, we consider the following amusing example taken from [2]. Let

$$f(\theta) = \sin^2(\theta) \sin^2(2\theta) \exp(-\theta^2/2) \quad (38)$$

be the (unnormalized) standard normal distribution perturbed by the factor $\sin^2(\theta) \sin^2(2\theta)$.

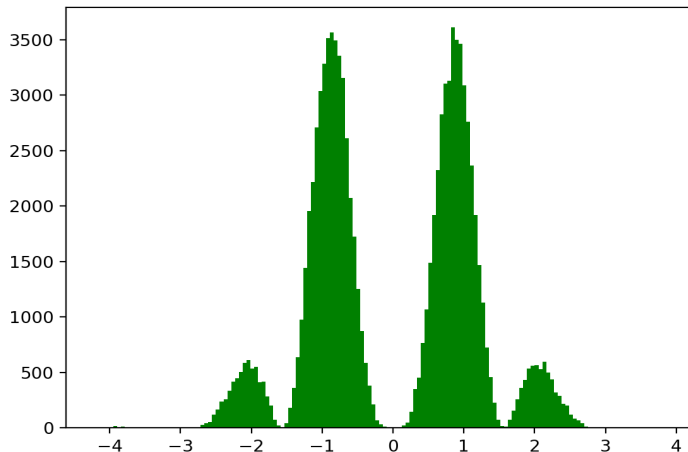
- We choose the proposal to be the uniform PDF,

$$q(\eta|\theta) = \frac{1}{2a} 1_{(\theta-a, \theta+a)}(\eta), \quad (39)$$

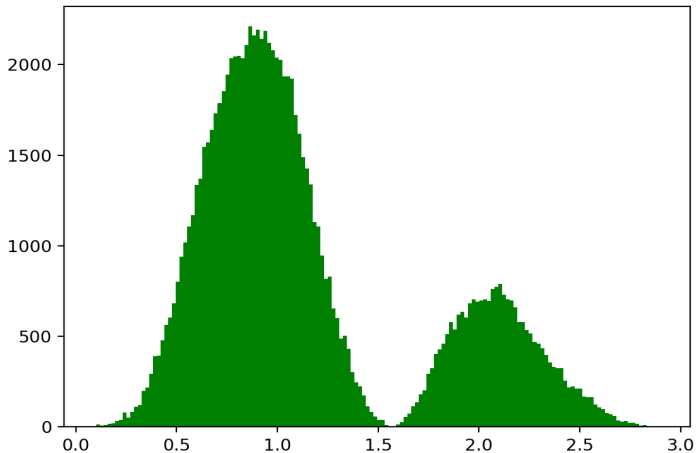
where $a > 0$ is a tunable parameter.

- The following graphs show the results of the simulations with two choices of the parameter a : $a = 1$ and $a = 0.1$, respectively. In both cases, $\theta_0 = 3.14$.
- Notice how the choice of a affects the simulation.

Example: perturbed Gaussian distribution



Example: perturbed Gaussian distribution



Example: perturbed Gaussian distribution

- Here is the code snippet used to generate these graphs in Python:

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




def target_density(x):
    fac = np.sin(x) * np.sin(2*x)
    return fac**2 * np.exp(-0.5*x**2)

def metropolis(theta, a):
    eta = np.random.uniform(theta - a, theta + a)
    u = np.random.uniform()
    if u > target_density(eta) / target_density(theta):
        eta = theta
    return eta

N = 100000
theta = np.zeros(N)
theta[0] = 3.14
for i in range(1, N):
    theta[i] = metropolis(theta[i-1], 1.0)

plt.hist(theta[1000:N], 150, color = 'g')
plt.show()
```

References

-  [1] *Handbook of Markov Chain Monte Carlo*, ed. by S. Brooks, G. L. Jones, and X.-L. Meng, CRC Press (2011).
-  [2] Robert, C. P.: The Metropolis-Hastings algorithm, arXiv:1504.01896v3 [stat.CO] 27 Jan 2016.
-  [3] Robert, C. P. and Casella, G.: *Monte Carlo Statistical Methods*, Springer (2004).