

Distribution	Bernoulli
Description	The random variable only has 2 possible outcomes. Probability of one of them is $p$
Notation	$X \sim \text{Bernoulli}(p)$
PMF	$P(X = k) = \begin{cases} p, & k = 1 \\ 1 - p, & k = 0 \end{cases}$
Expectation	$E(X) = p$
Variance	$\text{Var}(X) = p(1 - p)$
Properties	Indicator Function is a Bernoulli Random Variable, $1_A = \begin{cases} 1, & \text{if } A \text{ happens} \\ 0, & \text{if } A \text{ doesn't happen} \end{cases}$

Distribution	Binomial
Description	Number of successes in $n$ Bernoulli trials.
Notation	$X \sim B(n, p)$
PMF	$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n$
Expectation	$E(X) = np$
Variance	$\text{Var}(X) = np(1 - p)$
Properties	If $X_1, \dots, X_n$ are i.i.d. with distribution <i>Bernoulli</i> ( $p$ ), then $X_1 + \dots + X_n \sim B(n, p)$

Distribution	Geometric
Description	Number of Bernoulli trials to obtain the <b>first</b> success.
Notation	$X \sim \text{Geometric}(p)$
PMF	$P(X = k) = p(1 - p)^{k-1}, \quad k = 1, 2, 3, \dots$
Expectation	$E(X) = \frac{1}{p}$
Variance	$\text{Var}(X) = \frac{1-p}{p^2}$
Distribution	Poisson
Description	Number of events occurring in a fixed time interval or region of opportunity. Number of events per single unit of time.
Notation	$X \sim \text{Poi}(\lambda)$
PMF	$P(X = k) = \frac{e^{-\lambda}}{k!} e^{-\lambda}, \quad k = 0, 1, 2, \dots, \lambda > 0$
Expectation	$E(X) = \lambda$
Variance	$\text{Var}(X) = \lambda$
Properties	When $n$ is large & $p$ is small, $np$ is moderate $B(n, p) \rightarrow \text{Poi}(np)$

Distribution	Negative Binomial
Description	Number of Bernoulli trials to obtain $r$ successes.
Notation	$X \sim NB(r, p)$
PMF	$P(X = k) = \binom{r-1}{k-1} p^r (1 - p)^{k-r}, \quad k = r, r + 1, \dots$
Expectation	$E(X) = \frac{r}{p}$
Variance	$\text{Var}(X) = \frac{r(1-p)}{p^2}$
Distribution	Multinomial
Description	$n$ : Number of trials, $k$ : Number of mutually exclusive events.
Notation	$X \sim \text{Multinomial}(n, p_1, \dots, p_k)$
PMF	$P(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$
Alternate PMF	$p(\mathbf{x}) = \frac{n!}{\prod_j x_j!} \prod \pi_j^{x_j}$ for $\sum \pi_j = 1, \quad \sum x_j = n, x_j \geq 1$
Expectation	$E(X_i) = np_i$
Variance	$\text{Var}(X_i) = np_i(1 - p_i)$
Covariance	$\text{Cov}(X_i, X_j) = -np_i p_j, \quad i \neq j$

Distribution	Uniform
Notation	$X \sim \text{Uniform}(a, b)$
PDF	$f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$
CDF	$F(x) = \begin{cases} \frac{x-a}{b-a}, & a \leq x < b \\ 1, & x \geq b \end{cases}$
Expectation	$E(X) = \frac{a+b}{2}$
Variance	$\text{Var}(X) = \frac{(b-a)^2}{12}$
Properties	$U(0, 1) \equiv \text{Beta}(1, 1)$ Transform to <i>Uniform</i> ( $a, b$ ) from $U(0, 1)$ : $Y = (b - a)X + a$
Distribution	Normal
Notation	$X \sim N(\mu, \sigma^2)$
PDF	$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad -\infty < x < \infty$
CDF	$F(x) = \int_{-\infty}^x f(x)dx; -\infty < x < \infty$
Expectation	$E(X) = \mu$
Variance	$\text{Var}(X) = \sigma^2$
Properties	If $Z \sim N(0, 1)$ , then $\mu + \sigma Z \sim N(\mu, \sigma^2)$ Can complete the square if we have single exp then we

1.2 Integration Transformation:	$[0, \infty] \rightarrow [0, 1]$
Substitute $y = \frac{x}{x+1}, \quad dy = \frac{dx}{(x+1)^2} = -y^2 dx$	
$\theta = \int_0^\infty g(x)dx = \int_0^1 h(y)dy$	
$h(y) = \frac{g(\frac{y}{1-y})}{y^2}$	
1.3 Multi-dimensional Integral:	
$\theta = \int_0^1 \int_0^1 \dots \int_0^1 g(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n$	
Generate $k$ independent sets with $n$ independent uniform R.V. (i.e. $k \times n$ matrix)	

$\begin{pmatrix} u_{(1)}^{(1)} & \dots & u_{(1)}^{(1)} \\ u_{(1)}^{(2)} & \dots & u_{(2)}^{(2)} \\ \vdots & \dots & \vdots \\ u_{(1)}^{(k)} & \dots & u_{(k)}^{(k)} \end{pmatrix}$	
--	--

Since  $g(U_1^{(i)}, \dots, U_n^{(i)}), \quad i = 1, 2, \dots, k$  are i.i.d,  $\theta = \frac{1}{k} \sum_{i=1}^k g(U_1^{(i)}, \dots, U_n^{(i)}) \rightarrow \theta$

#### 1.4 Algorithm for estimating $\pi$

- Suppose  $(X, Y)$  is uniformly distributed in the square of area 4 centred at the origin (0,0).
- $P(X^2 + Y^2 \leq 1) = \frac{\pi}{4}$
- Hence, if a large number of points are generated in the square, the proportion of points that fall within the circle is approximately  $\frac{\pi}{4}$

#### 2. Inversion Method Algorithm

- Generate a random number  $u = U \sim [0, 1]$
- If R.V. is continuous, find  $X$  in which  $F_X^{-1}(u)$  falls under
- Else, for  $i = 0, 1, 2, \dots, n$ , find  $X$  where  $F_X = \sum_{i=0}^{i=n} p_i$
- Set  $X = F_X^{-1}(u)$

#### 2.1 Proof for Inversion Method

- Let  $F_X$  denote the distribution of  $X = F^{-1}(U)$ .
- $F_X(x) = P(X \leq x) = P(F^{-1}(U) \leq x)$ .
- Since  $F$  is a distribution function,  $F(x)$  is a monotone increasing function of  $x$ . So  $a \leq b \implies F(a) \leq F(b)$
- Therefore,  $F_X(x) = P[F^{-1}(U)] \leq F(x)) = P(U \leq F(x)) = F(x)$ . Since  $U \sim [0, 1]$ 
  - Use when theoretical  $F^{-1}$  is known
  - If discrete, make sure to calculate the individual  $F_X$  properly.
  - At times, calculating all the individual  $F_X$  is difficult, like with Poisson, Binomial & Negative Binomial
  - For Gamma, Inversion method can only be used if you convert the Gamma to a sum of independent Negative Binomials.
  - In continuous case, use  $\int_a^b f(y)dy$  or cdf to calculate.

#### 2.2 Inversion Method Table

For  $S \sim \text{Gamma}(n, \beta)$ , let  $S_n = X_1 + X_2 + \dots + X_n$ , where  $X_i \sim \text{Exp}(\beta)$

Name	Base Case	Recursion Relation
Poisson	$P(X = 0) = e^{-\lambda}$	$P(X = k) = P(X = k-1) \cdot \frac{\lambda}{k}, k = 1, 2, \dots$
Binomial	$P(X = 0) = (1 - p)^n$	$P(X = k) = P(X = k-1) \cdot \frac{(n-k+1)p}{k(1-p)}, \quad k = 1, 2, \dots, n$
Negative Binomial	$P(X = r) = p^r$	$P(X = k) = P(X = k-1) \cdot \frac{k-k+r}{k-k+r-1} \cdot (1-p)$

#### 3. Rejection Method Algorithm

- Choose a **simpler** probability density function  $g(x)$  (*Trial/Proposal pdf*), such that  $f_X(x) \leq c \cdot g(x), \forall x$  in  $x : f_X(x)$  where  $c > 1$ .
- Generate a random variable  $V$  from  $g(x)$
- Generate a uniform random number  $Y$  from  $U \sim [0, c \cdot g(x)]$ .
- If  $Y \leq f_X(V)$ , accept  $V = X$ .
- Else, reject  $V$  and repeat generation of  $V$  &  $Y$  until  $Y \leq f_X(V)$  holds.

#### 3.1 Proof for Rejection Method

- $(V, Y)$  is uniformly distributed over  $S = \{(x, y) : 0 \leq y \leq c \cdot g(x)\}$
- $(V, Y)$  is only accepted if it lies in  $B = \{(x, y) : 0 \leq y \leq f(x)\}$
- Since  $(V, Y)$  is uniform over  $S$ , the accepted  $(V, Y)$  is uniformly distributed over  $B$
- The joint density of  $(V, Y) : h(x, y) = \begin{cases} g(x)h(y|x) = \frac{1}{c}, & (x, y) \in S \\ 0 & \end{cases}$ 
  - The chances of sampling  $(V, Y)$  from  $S$  is  $\frac{1}{c}$
- The marginal density of accepted  $V : k(x) = \int_0^{f(x)} 1 dy = f(x)$

$\cos \theta = \frac{V_x}{R} = \frac{V_x}{\sqrt{V_x^2 + V_y^2}}$
$\sin \theta = \frac{V_y}{R} = \frac{V_y}{\sqrt{V_x^2 + V_y^2}}$
Following Box-Muller
$X = \sqrt{-2\log(U)} \cdot \frac{V_1}{\sqrt{V_1^2 + V_2^2}}$
$Y = \sqrt{-2\log(U)} \cdot \frac{V_2}{\sqrt{V_1^2 + V_2^2}}$

Given  $R = X^2 + Y^2, R \sim U[0, 1]$  and is independent from  $\theta$ , let  $R = S$ ,

$X = \sqrt{\frac{-2\log(U)}{V_1^2 + V_2^2}} \cdot V_1 = \sqrt{\frac{-2\log(S)}{S}} V_1$
$Y = \sqrt{\frac{-2\log(U)}{V_1^2 + V_2^2}} \cdot V_2 = \sqrt{\frac{-2\log(S)}{S}} V_2$

$X$  and  $Y$  are independent unit normals when  $(V_1, V_2)$  is a randomly chosen point in the circle of radius 1 centred at the origin and  $S = V_1^2 + V_2^2$

**5.5 Multivariate Normal R.V. Algorithm**  
 $\Sigma$  is a  $d \times d$  covariance matrix

- $X \sim N_d(\mu, \Sigma)$  has p.d.f:  $p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} e^{-\frac{(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}{2}}$

- Generate vector  $\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{pmatrix}$  where  $\mathbf{z} \in \mathbb{R}^d \quad z_i \sim N(0, 1)$

- $\mathbf{z} \sim N(0, I)$  where  $I$  is identity matrix
- Find a  $d \times d$  matrix  $T$  where  $T^T T = \Sigma$
- Let  $\mathbf{x} = T^T \mathbf{z} + \mu$ 
  - Matrix  $T$  can be found by Cholesky decomposition
  - Matrix  $T$  can be found by Eigenvalue decomposition
  - This does not use Uniform number generator but `rnorm` generators

#### 5.6 Multivariate Conditional Normal R.V. Algorithm

- $X \sim N_d(\mu, \Sigma)$  has p.d.f:  $p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} e^{-\frac{(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}{2}}$
- Generate  $x_1$  from  $X_1 \sim N(\mu_1, \sigma_{11})$
- Generate  $x_2|x_1$  from  $X_2|X_1 \sim N(\mu_2', \sigma_{22}')$
- Generate  $x_3|(x_2, x_1)$  from  $X_3|(X_1, X_2) \sim N(\mu_3', \sigma_{33}')$
- Repeat until the full vector  $\mathbf{x}$  is solved
- $\mu_i'$  is the conditional mean of  $x_i$  represented as a vector
- $\Sigma_{ii}'$  is the conditional variance of  $x_i$  represented as a vector
- $\mu_i' = \mu_{i|i-1} = \mu_i + \Sigma_{i,1:i-1} \Sigma_{1:i-1,1:i-1}^{-1} (\mathbf{x} - \mu_{1:i-1})$
- $\Sigma_{ii}' = \Sigma_{ii|i-1} = \Sigma_{ii} - \Sigma_{i,1:i-1} \Sigma_{1:i-1,1:i-1}^{-1} \Sigma_{1:i-1,i}$
- This method uses the fact that multivariate normals have a joint Gaussian relationship (i.e. conditional in nature)

#### 5.7 Multinomial Algorithm

- Multinomial :  $p(\mathbf{x}) = \frac{n!}{\prod_j x_j!} \prod \pi_j^{x_j}$  for  $\sum \pi_j = 1, \quad \sum x_j = n, x_j \geq 1$
- Alternate:  $p(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! \dots x_k!} \cdot (p_1)^{x_1} \dots (p_k)^{x_k}$
- Generate  $x_1$  from  $X_1 \sim \text{Binomial}(n, p_1)$
- Generate  $x_j|x_1, \dots, x_{j-1}$  from  $X_j|X_1, \dots, X_{j-1} \sim \text{Binomial}(n - \sum_{j=1}^{j-1} X_i, \frac{p_j}{1 - \sum_{i=1}^{j-1} p_i})$
- Assign last  $x_k = n - \sum_{k=1}^{k-1} X_i$
- This method makes use of the fact that multinomial are essentially a sum of marginal binomials

#### 6. Variance Reduction

Let  $h(x)$  be any function used to represent the Expectation. For instance, if I am finding  $E[X]$  then  $h(X) = X$ . The goal is in reducing the value of  $n$  in *Monte Carlo Integration*, where the distribution  $X$  has the samples  $\mathbf{x} = [x_1, \dots, x_n]$

$$E[X] = E[h(\mathbf{x})] = \frac{1}{n} \int h(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = \theta$$

$$E[X^2] = E[h(\mathbf{x})^2] = \frac{1}{n} \int h(\mathbf{x})^2 f(\mathbf{x}) d\mathbf{x}$$

$$\text{Var}[X] = \text{Var}[h(\mathbf{x})] \approx \frac{E[(h(\mathbf{x}) - E[h(\mathbf{x})])^2]}{n} = \frac{E[h(\mathbf{x})^2] - (E[h(\mathbf{x})])^2}{n}$$

for some limits ("boundary") where area under the curve  $f(x) = 1$  as  $n$  approaches the confidence interval of the estimator.

Let  $\tilde{\mathbf{I}}$  be an MC estimator for some distribution.

#### 6.1 Simple Sampling Algorithm

Sample  $n$  samples from distribution  $X$

- $\tilde{\mathbf{I}} = \frac{1}{n} \sum_{i=1}^n h(X)$
- $\mathbf{I} = E[\tilde{\mathbf{I}}] = \tilde{\mathbf{I}}$
- $E[\tilde{\mathbf{I}}] = \frac{1}{n} \int h(X)^2 f(x) dx$
- $\text{Var}[\tilde{\mathbf{I}}] = \frac{\text{Var}[\mathbf{I}]}{n} = \frac{E[\tilde{\mathbf{I}}^2] - \tilde{\mathbf{I}}^2}{n}$
- Does not reduce variance
- Not an unbiased estimate of variance

### 7. Markov Chain

For our purposes, we assume that Markov Chains are homogenous in time, meaning the transition probabilities are fixed and do not change with time. A **Markov Chain** is a mathematical model used to describe a sequence of events where,

- The future state depends **only** on the current state (not the past). This is called the **Markov Property**

A Markov Chain  $X$  is a discrete time stochastic process  $X_0, X_1, \dots$  with the property that the distribution of  $X_t$  given all previous values of the process,  $X_0, \dots, X_{t-1}$ , only depends upon  $X_{t-1}$ :  $P[X_t \in A|X_0, \dots, X_{t-1}] = P[X_t \in A|X_{t-1}]$  for any Set  $A$

- The transitions between states are described by **probabilities**  $p_{ij} = P(X_{t+1} = j|X_t = i)$  denotes the transition probability from state  $i$  to state  $j$  at time  $t + 1$ .

$p_{ij}(m) = P(X_{t+m} = j|X_t = i), m = 1, 2, \dots$  denote the multi-step transition probability from state  $i$  to state  $j$ . Can be seen as the sum over all intermediate states  $k$  through which the system passes in its transition from state  $i$  to state  $j$

Recursive relation:

$$p_{ij}(m + 1) = \sum_k^n p_{ik}(m) p_{kj}, \quad m = 1, 2, \dots \text{ with } p_{ik}(1) = p_{ik}$$

$$p_{ij}(m + n) = \sum_k^n p_{ik}(m) p_{kj}(n), \quad m, n = 1, 2, \dots$$

#### 7.1 Markov Chain Components

- States:** A Markov Chain consists of a finite set of states,  $S = \{s_1, s_2, \dots, s_n\}$
- Initial State:** The starting state of the system is represented by a probability vector  $\pi^{(0)}$ , where  $\pi_i^{(0)}$  is the probability of starting in state  $i$
- Stationary Distribution/Invariant:** A probability distribution  $\pi$  is **stationary** if, after applying the transition matrix  $P$ , it remains unchanged:  $\pi P = \pi$ . Intuitively, it describes the long-term behaviour of the system.
- Transition Matrix:** The transitions between states are represented by a square matrix  $P$ , called the **transition matrix**. Each entry  $p_{ij}$  gives the probability of transitioning from state  $i$  to state  $j$ . The matrix is written as,

$$P = \begin{bmatrix} p_{11} & p_{12} & \dots & p_{1n} \\ p_{21} & p_{22} & \dots & p_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ p_{n1} & p_{n2} & \dots & p_{nn} \end{bmatrix}$$

All  $p_{ij}$  are conditional probabilities that satisfy  $p_{ij} \geq 0$  and  $\sum_{j=1}^n p_{ij} = 1$  where  $n$  is the total number of states.  $p_{ij} \geq 0$ : Probabilities are non-negative.  $\sum_{j=1}^n p_{ij} = 1$ : Each row sums to 1, ensuring the probabilities account for all possible transitions.

#### 7.3 Markov Chain Mechanics

- Current State:** At any time  $t$ , the system is in one of the states  $s_i \in S$ .
- Transition:** The system moves to a new state at the next time step  $t + 1$ , based on the probabilities in the row of  $P$  corresponding to the current state.
- Evolution:** The state probabilities at time  $t + 1$ , denoted  $\pi^{(t+1)}$ , are given by:  $\pi^{(t+1)} = \pi^{(t)} P$

$\pi$  is a vector

#### 7.4 Markov Chain Terminologies

- Irreducible:** A Markov Chain is irreducible if it's possible to get from any state to any other state (possibly over multiple steps).
- Aperiodic:** A chain is aperiodic if it doesn't repeat states in a fixed cycle (e.g., alternating between two states indefinitely).
- Recurent States:** States that the system will return to infinitely often.  $f_i = P(\text{ever returning to state } i) = 1$
- Transition States:**  $f_i < 1$
- Transient States:** States that the system may leave and never return to.
- Accessible States:** if  $i$  can reach  $j$  in a finite number of transitions or vice versa
- Communicable States:**  $i \leftrightarrow j$  or if  $i$  and  $j$  can transition to each other in a finite number of transitions or vice versa
- Positive Recurrent:** A recurrent state is positive recurrent if the expected time to return to that state is finite.

#### 7.5 Markov Chain Example

Suppose we have three states for the weather:  $S = \{\text{Sunny, Rainy, Cloudy}\}$

The transition matrix  $P$  might look like:

$$P = \begin{bmatrix} 0.6 & 0.3 & 0.1 \\ 0.4 & 0.4 & 0.2 \\ 0.3 & 0.3 & 0.4 \end{bmatrix}$$

If it's **Sunny** today,

### 10. Simulated Annealing Algorithm

This seeks to find a global min/max of a cost function by exploring the solution space whilst avoiding getting trapped in local minima/maxima

- Set a sufficiently high initial temperature  $T_0$  to ensure most proposed transitions are accepted, regardless of cost function improvement.
- Reduce temperature gradually following an exponential cooling schedule:  $T_k = \alpha T_{k-1}, \quad k = 1, 2, \dots$  where  $T_k$ : Temperature at step  $k$ ,  $\alpha$  is a constant reduction factor(0.8 to 0.99). Smaller values  $\rightarrow$  faster cooling, values close to  $1 \rightarrow$  slower cooling.
- Initially, large changes are accepted, but as the temperature decreases, the algorithm becomes more selective.
- At each temperature  $T_k$ , a number of proposed transitions (new solutions) are attempted.
- If the new solution improves the cost function, it is always accepted.
- If the new solution worsens the cost function, it is accepted with a probability proportional to  $\exp(-\Delta E/T)$ , where  $\Delta E$  is the change in the cost function. This helps to escape local minima early in the process.
- Continue until no significant changes in the solution or the desired number of acceptances is not achieved for several consecutive temperature levels (Convergence)

**11. Confidence Intervals**  
*Estimate Mean:*  $\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N X_i$ , where  $X_i$  are the simulated values and  $N$  is the number of samples.

*Estimate Variance:*  $S^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \hat{\mu}_N)^2$

*Compute Standard Error:*  $SE = \sqrt{\frac{S^2}{N}}$   
*Determine Critical Value:* For a 95% confidence interval,  $z^* \approx 1.96$  (standard normal distribution).

*Confidence Interval:* 95% confidence interval is given by:  $\hat{\mu}_N \pm z^* \cdot SE$ . This translates to the interval:  $[\hat{\mu}_N - z^* \cdot SE, \hat{\mu}_N + z^* \cdot SE]$   
*Adjust for Dependencies (if not i.i.d.):* For stationary sequences, adjust the variance using the effective sample size:  $\text{Var}(\hat{\mu}_N) \approx \frac{\text{Var}(X_i) \cdot \tau}{N}$ , where  $\tau$  is the integrated autocorrelation time.

#### 11.1 Confidence Intervals for Case 2: Batch Means

When dealing with stationary processes, consecutive observations are **correlated**, making the usual formula for variance unreliable.

- Split  $N$  total observations  $X_1, \dots, X_N$  into  $b$  batches of size  $L = N/b$
- Compute the mean for each batch:  $\bar{Y}_k = \frac{1}{L} \sum_{i=(k-1)L+1}^{kL} X_i$
- Treat these batch means  $\bar{Y}_1, \dots, \bar{Y}_b$  as i.i.d.
- Estimate the variance of  $\bar{X}$  (mean of the entire sequence) using the variance of  $\bar{Y}_k$ :  $S_b^2 = \frac{1}{b-1} \sum_{k=1}^b (Y_k - \bar{Y}_b)^2$
- Construct a 95% confidence interval:  $\bar{Y}_b \pm t_{b-1, 0.025} \sqrt{\frac{S_b^2}{b}}$  where  $t_{b-1, 0.025}$  is the critical value of the  $t - \text{distribution}$  with  $b - 1$  degrees of freedom.
- The batch size  $L$  must be large enough to ensure the batch means  $\bar{Y}_k$  are approximately independent.

#### 11.2 Confidence Intervals for Case 2: Covariance Summation

For a stationary process, the variance of the sample mean  $\bar{X}_N$  depends on the **autocovariance** structure of the process:

$$\text{Var}(\bar{X}_N) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \text{Cov}(X_i, X_j)$$

- Define  $c(k) = \text{Cov}(X_1, X_{1+k})$
- For stationary processes:  $\text{Var}(\bar{X}_N) = \frac{1}{N} \sum_{k=-N}^{N-1} \left(1 - \frac{|k|}{N}\right) c(k)$
- For large  $N$ , this simplifies to:  $\text{Var}(\bar{X}_N) \approx \frac{1}{N} V$  where  $V = c(0) + 2 \sum_{k=1}^\infty c(k)$
- Directly estimating  $V$  using:  $\hat{V}_N = \hat{c}_N(0) + 2 \sum_{k=1}^{N-1} \hat{c}_N(k)$  where  $c^N(k) \hat{c}_N(k)$  is the sample covariance at lag  $k$ , often performs poorly in practice due to noise and bias.

#### 11.3 Confidence Intervals for Case 2: Window Methods

Instead of summing the entire auto-covariance series  $V = c(0) + 2 \sum_{k=1}^\infty c(k)$  (which is computationally expensive and may not converge well), you sum up to a fixed lag  $L$  (called the "window size"):  $\hat{V}_{N,L} = \hat{c}_N(0) + 2 \sum_{k=1}^L \hat{c$

	see the distribution of it and the normalising constant is straightforward
<b>Distribution</b>	Exponential
<i>Notation</i>	$X \sim \text{Exp}(\lambda)$
<i>PDF</i>	$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$ <p>Note that <math>\lambda &gt; 0</math></p>
<i>CDF</i>	$F(x) = \begin{cases} 1 - e^{-\lambda x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$
<i>Expectation</i>	$E(X) = \frac{1}{\lambda}$
<i>Variance</i>	$\text{Var}(X) = \frac{1}{\lambda^2}$
<i>Properties</i>	For any $X \sim \text{Exp}(\lambda)$ $P(X > s + t \mid X > s) = P(X > t)$
<b>Distribution</b>	Gamma
<i>Notation</i>	$X \sim \text{Gamma}(\alpha, \beta)$
<i>PDF</i>	$f(x) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
<i>CDF</i>	$F(x) = \frac{1}{\Gamma(\alpha)} \gamma(\alpha, \beta x)$
<i>Expectation</i>	$E(X) = \frac{\alpha}{\beta}$
<i>Variance</i>	$\text{Var}(X) = \frac{\alpha}{\beta^2}$
	$\gamma(\alpha, \beta x) = \int_0^{\beta x} t^{\alpha-1} e^{-t} dt$ only computable where we know $x$
	<b>Gamma Function:</b> $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$
	<b>Recursion Property:</b> $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$
	<b>Gamma Function Computation:</b> $\Gamma(\alpha) = (\alpha - 1)!$
	<b>Sum of Gamma Random Variables:</b> Let $X_1, \dots, X_k$ be independent random variables where $X_i \sim \text{Gamma}(\alpha_i, \beta)$ for each $i$ . Then $X_1 + \dots + X_k \sim \text{Gamma}(\alpha_1 + \dots + \alpha_k, \beta)$
<i>Properties</i>	<b>Connection with the Standard Normal:</b> If $Z \sim N(0, 1)$ , then: $Z^2 \sim \text{Gamma}(\frac{1}{2}, \frac{1}{2}) \sim \chi^2(1)$ <b>Connection with Chi Squared:</b> Assume $Z_1, \dots, Z_k$ are i.i.d. $N(0, 1)$ random variables. Then, $Z_1^2 + \dots + Z_k^2 \sim \text{Gamma}(\frac{k}{2}, \frac{1}{2}) \sim \chi^2(k)$ <b>Connection with Exponential:</b> If $X_1, \dots, X_n$ i.i.d. $\text{Exp}(\lambda) = \text{Gamma}(1, \lambda)$ then $X_1 + \dots + X_n \sim \text{Gamma}(n, \lambda)$ <b>Scaling:</b> If $X \sim \text{Gamma}(\alpha, \beta)$ , then $cX \sim \text{Gamma}(\alpha, \frac{\beta}{c})$

<b>Distribution</b>	Beta
<i>Notation</i>	$X \sim \text{Beta}(\alpha, \beta)$
<i>PDF</i>	$f(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$ for $0 \leq x \leq 1, \alpha > 0, \beta > 0$
<i>CDF</i>	$F(x) = \frac{1}{\Gamma(\alpha)} \gamma(\alpha, \beta x)$
<i>Expectation</i>	$E(X) = \frac{\alpha}{\alpha+\beta}$
<i>Variance</i>	$\text{Var}(X) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
	<b>Swap of Parameters:</b> If $X \sim \text{Beta}(\alpha, \beta)$ , then $1 - X \sim \text{Beta}(\beta, \alpha)$
	<b>Connection with Gamma Distribution:</b> If $X \sim \text{Gamma}(\alpha, \beta)$ , $Y \sim \text{Gamma}(\beta, \beta)$ , and $X, Y$ are independent, then $\frac{X}{X+Y} \sim \text{Beta}(\alpha, \beta)$
<i>Properties</i>	<b>Order Statistics:</b> If $X_1, \dots, X_n$ are i.i.d. $\text{Uniform}(0, 1)$ , and $X_{(1)} \leq \dots \leq X_{(n)}$ are their order statistics, then for $k = 1, \dots, n$ , $X_{(k)} \sim \text{Beta}(k, n + 1 - k)$ <b>Useful Tip:</b> To generate Beta distribution, draw i.i.d. Uniform samples, order them, and select the $k$ th one.

#### 1. Monte Carlo Integration

- Suppose  $\theta = E[g(x)] = \int_0^1 g(x) dx$
- Assuming  $U \sim [0, 1]$ ,  $\theta = E[g(U)]$
- If  $U_1, U_2, \dots, U_k$  are independent  $U$ ,  $g(U_1), g(U_2), \dots, g(U_k)$  are i.i.d
- By **Strong Law of Large Numbers**,  $\frac{1}{k} \sum_{i=0}^k g(U_i) \rightarrow E[g(U)] = \theta$
- Generate a large number of random numbers  $U_i$
- Approximate  $\theta$  by the average value of  $g(U)$  with  $\hat{\theta} = \frac{1}{k} \sum_{i=0}^k g(U_i)$

#### 1.1 Integration Transformation: $[a, b] \rightarrow [0, 1]$

Substitute  $y = \frac{x-a}{b-a}$ ,  $dy = \frac{dx}{b-a}$   
 $\theta = \int_a^b g(x) dx = \int_0^1 g(a + [b-a]y) (b-a) dy = \int_0^1 h(y) dy$

- $B$  is contained by  $f(x)$ , it's total area under the curve is 1 as  $f(x)$  is some pdf
- Integrating  $B$  over the region bounded by  $f(x)$  (i.e.  $[0, f(x)]$ ) will be where  $V$  is found
- Uniform sampling ensures proportionality between  $f(x)$  and the region  $B$
- Efficient when  $g$  is chosen s.t.  $c$  is small - close to 1 (fewer points are rejected)
- $c$  should be the maximum value of  $\frac{f(x)}{g(x)}$  over the entire domain of  $x$  (i.e.  $c = \sup \frac{f(x)}{g(x)}$  or use  $\frac{d}{dx} = 0$ )
- Number of proposals needed is  $c$  where  $c \sim \text{Gemoetric}(\frac{1}{c})$
- $P(\{(V, Y) : B\}) = \frac{\int_{\text{area } B} c}{\int_{\text{area } S} c} = \frac{1}{c}$

#### 4. Composition Method Algorithm

- Divide  $F_X(x)$  into  $M$  subregions where  $F_X(x) = \sum_{i=1}^M p_i \cdot F_{X_i}(x)$  where  $\sum_{i=1}^M p_i = 1$
- Choose an  $F_{X_i}$ , where  $i \sim$  some discrete r.v. (use uniform for most cases but depends on distribution of  $\sum_{i=1}^M p_i = 1$ )
- Generate  $Y = U \sim [0, 1]$  and scale it to the range of  $F_{X_i}$
- Solve  $X = F_{X_i}^{-1}(Y)$

#### 4.1 Composition Method Remarks

- Useful if you know how  $F_X^{-1}$  can be "binned" into weighted distributions
- Useful when  $F_{X_i}^{-1}$  is known but not  $F_X^{-1}$

#### 5. Polar Method Algorithm for Generating Normal R.V.

Using the relationship between cartesian  $(x, y)$  and polar  $(r, \theta)$  coordinates, Normal R.V. can be simulated as follows

##### 5.1 Box-Muller Algorithm

- Generate  $U_1, U_2 \sim U[0, 1]$
- Since  $R \sim \text{Exp}(\frac{1}{2})$ ,  $R = F^{-1}(U_1) = -2\log(U_1)$
- Since  $\Theta \sim U[0, 2\pi]$ ,  $\theta = F^{-1}(U_2) = 2\pi * U_2$
- Let  $X = R \cos \theta = \sqrt{-2\log(U_1)} \cos 2\pi * U_2$
- Let  $Y = R \sin \theta = \sqrt{-2\log(U_1)} \sin 2\pi * U_2$

##### 5.2 Box-Muller Algorithm Proof

Suppose  $X, Y \sim N(0, 1)$  where  $X$  and  $Y$  are i.i.d.  
Let  $R$  &  $\Theta$  be polar coordinates of vector  $(X, Y)$ .  
 $R = \sqrt{X^2 + Y^2}$ ,  $\tan \Theta = \frac{Y}{X}$   
Given  $X$  and  $Y$  are independent, joint density is,  
 $f(x, y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} = \frac{1}{2\pi} e^{-\frac{(x^2+y^2)}{2}}$   
Convert the equations to a "realisation" format, with Jacobian Transformation of  $|J| = 2$  accounting for the change in  $(x, y)$  coordinates to  $(r, \theta)$  coordinates.  
 $r = x^2 + y^2$ ,  $\theta = \tan^{-1} \frac{y}{x}$   
 $f(r, \theta) = \frac{1}{2\pi} e^{-\frac{r}{2}}$ ,  $0 < r < \infty$ ,  $0 < \theta < 2\pi$   
The joint density of  $(r, \theta)$  can be further factored into 2 parts, an exponential and a uniform density,  $f(r, \theta) = \frac{1}{2} \frac{1}{2\pi} e^{-\frac{r}{2}} = (\frac{1}{2} e^{-\frac{r}{2}}) \cdot (\frac{1}{2\pi})$   
This means  $R \sim \text{Exp}(\frac{1}{2})$  and  $\Theta \sim U[0, 2\pi]$

Using this, we can generate  $(r, \theta)$  coordinates and convert it into  $(x, y)$  coordinates

##### 5.3 Improved Box-Muller Algorithm

- Generate  $U_1, U_2 \sim U[0, 1]$
- Let  $V_1 = 2U_1 - 1$
- Let  $V_2 = 2U_2 - 1$
- Let  $S = U_1^2 + V_2^2$
- If  $S > 1$ , repeat
- Else,  $X = \sqrt{\frac{-2\log(S)}{S}} V_1$  and  $Y = \sqrt{\frac{-2\log(S)}{S}} V_2$

##### 5.4 Improved Box-Muller Algorithm Proof

Given Step 1 - 3,  $(V_1, V_2)$  will be uniformly distributed in a square centred on  $(0, 0)$  with area 4. Suppose,  $(V_1, V_2)$  continue to be generated until it is in a circle centred on  $(0, 0)$  of radius 1 or  $V_1^2 + V_2^2 \leq 1$ . Hence,  $(V_1, V_2)$  following this condition is uniformly distributed in the circle.  
If  $(V_1, V_2)$  are converted to polar coordinates using Jacobian Transformation,  
 $f(r, \theta) = (1) \cdot (\frac{1}{2\pi})$   
This means  $R \sim U[0, 1]$  and  $\Theta \sim U[0, 2\pi]$  where  $R$  and  $\Theta$  are independent

Since  $\theta$  is a random angle,  $(V_1, V_2)$  being points randomly generated in the circle, the  $\sin \theta$  and  $\cos \theta$  can be generated

#### 6.2 Stratified Sampling Algorithm

- Divide domain of  $f(x) : S$  into  $M$  disjoint subsets  $S^{(1)}, S^{(2)}, \dots, S^{(M)} : S = \bigcup_{i=1}^M S^{(i)}$ .
- For each  $S^{(i)} : P(X \in S^{(i)}) = a_i = \int_{S^{(i)}} f(x) dx$
- Assign  $n_i$  as sample sizes to each  $S^{(i)} : n_1 + \dots + n_M = n$
- For each  $S^{(i)}$ , generate  $n_i$  realisations  $X_1^{(i)}, \dots, X_{n_i}^{(i)}$  from

conditional PDF:  $g(x) = \begin{cases} \frac{f(x)}{a_i}, & \text{if } x \in S^{(i)}, \\ 0, & \text{otherwise.} \end{cases}$

- The sub-strata estimator is  $\hat{T}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} h(X_j^{(i)})$  for some function  $h$
- $E[T_i] = \int_{S_i} h(\mathbf{x}_{S_i}) \frac{f(\mathbf{x}_{S_i})}{a_i} = \frac{1}{a_i} \int_{S_i} h(\mathbf{x}_{S_i}) f(\mathbf{x}_{S_i}) = \frac{a_i}{a_i} = 1$
- The stratified estimator is  $T = \sum_{i=1}^M a_i T_i$  after correcting for sub-dividing  $n$  and drawing samples from  $g(x)$
- $E[T] = \sum_{i=1}^M a_i E[T_i] = \sum_{i=1}^M a_i \frac{a_i}{a_i} = I$
- $\text{Var}(T_i) = \frac{1}{n_i} \left( \int_{S^{(i)}} h(x)^2 \frac{f(x)}{a_i} dx - \left( \frac{a_i}{a_i} \right)^2 \right)$
- $\text{Var}T = \sum_{i=1}^M a_i^2 \text{Var}T_i$

- $a_i = \int_{S^{(i)}} h(x) f(x) dx$  is the total probability of  $S^{(i)}$  under  $f(x)$ .
- Dividing  $f(x)$  by  $a_i$  ensures the conditional PDF  $g(x)$  sums to 1 over  $S^{(i)} : \int_{S^{(i)}} g(x) dx = \int_{S^{(i)}} \frac{f(x)}{a_i} dx = \frac{1}{a_i} \int_{S^{(i)}} f(x) dx = 1$ .
- $T$  is unbiased because the sum of the sub-strat aestimators result in the estimator itself
- Reduces variance because  $\text{Var}\hat{\theta} = \text{Var}T + \frac{1}{n} \sum_{i=1}^M a_i (\frac{1}{a_i} - I_B)$  which is  $\geq \text{Var}T$

#### 6.1 Importance Sampling Algorithm

We are first given  $I = E[h(X)] = \int_S h(x) f(x) dx$  we can choose some  $g(x)$  and sample  $X_i$ s from  $g$ . It needs to be noted that  $g$  should be as close in shape to  $h(x)f(x)$  for this to work.  $Y$  is just to note that the samples do not follow the original distribution of  $X$  since we are drawing from  $g$ .

- Rewrite  $I = \int_S h(x) \frac{f(x)}{g(x)} g(x) dx = E[\frac{h(x)f(x)}{g(x)}] = E[h(Y)w(Y)]$  where  $w(y) = \frac{f(y)}{g(y)}$  is the weighting function.
- $\hat{I} = \frac{1}{n} \sum_{i=1}^n h(x_i) w(x_i)$  where  $w(x_i) = \frac{f(x_i)}{g(x_i)}$
- $\text{Var}\hat{I} = \frac{1}{n} \text{Var}[h(X)w(X)] = \frac{1}{n} (\int_S \frac{h^2(x)f^2(x)}{g(x)} dx - I^2)$ .
- If  $h(x) \geq 0$  for all  $x \in S$ , the  $\text{Var}\hat{I}$  is exactly 0.
- $\sigma^2 = \text{Var}g[h(X)w(X)] = \int_S \frac{h^2(x)f^2(x)}{g(x)} dx - I^2$
- Confidence Interval :  $I \in [\hat{I} - 1.96 \frac{\hat{\sigma}}{\sqrt{n}}, \hat{I} + 1.96 \frac{\hat{\sigma}}{\sqrt{n}}]$  where  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \frac{h^2(x_i)f^2(x_i)}{g(x_i)} - \hat{I}^2$
- Define  $g'(x) = c \cdot [h(x)] \cdot f(x)$  where  $c$  is a constant
- Let  $\int_S g'(x) dx = 1$ , hence  $c = \frac{1}{\int_S h(x)f(x) dx}$  or  $\frac{1}{I}$  if the estimator is known
- $g(x) = c \cdot g'(x)$
- Optimal proposal density is  $g(x) \propto |h(x)| \cdot f(x)$

#### 6.4 Control Variates Sampling Algorithm

We are first given  $I = E[h(X)] = \int_S h(x) f(x) dx$ . The control variates method uses an auxiliary variable  $C$ , correlated with  $X$ , to reduce the variance of the estimator

- Define the adjusted estimator as  $\hat{I} = \hat{I}_X + \beta(\hat{C} - \mu_C)$
- $\hat{I}_X = \frac{1}{n} \sum_{i=1}^n h(X_i)$
- $\hat{C} = \frac{1}{n} \sum_{i=1}^n C_i$
- $\mu_C = E[C]$  is known or approximated
- $\beta$  is a coefficient determined by minimising variance.
- Use optimal  $\beta^* = -\frac{\text{Cov}(h(X), C)}{\text{Var}(C)}$
- $\text{Var}(\hat{I}) = \text{Var}(\hat{I}_X) - 2\beta \text{Cov}(h(X), C) + \beta^2 \text{Var}(C) = (1 - \text{Cor}(h(X), C)^2) \cdot \text{Var}(\hat{I}_X)$
- If  $C$  is highly correlated with  $h(X)$ , a  $\beta$  can always be chosen so that  $\text{Var}(\hat{I}) \leq \text{Var}(\hat{I}_X)$
- if only  $E[C] = \mu$  is given,  $\hat{I} = \beta \hat{I}_X + (1 - \beta) C$

#### 6.5 Antithetic Variates Method

This method exploits negatively correlated random variables to reduce variance, particularly if  $h(X)$  is monotonic

- Let  $U \sim U[0, 1]$ ,  $X = F^{-1}(U)$  and  $X' = F^{-1}(1 - U)$
- $\hat{I} = \frac{1}{2n} \sum_{i=1}^n (h(U_i) + h(1 - U_i))$
- $\text{Var}[\hat{I}] = \frac{1}{2n} (\text{Var}(h(U)) + \text{Cov}(h(U), h(1 - U)))$ .
- For  $U \sim U[0, 1]$ ,  $U$  and  $1 - U$  are perfectly negatively correlated, reducing variance compared to independent sampling.
- The effectiveness of the antithetic method depends on the monotonicity of  $h(x)$ . If  $h(x)$  is monotonic,  $h(U)$  and  $h(1 - U)$  are negatively correlated, ensuring variance reduction

- 60% chance it will be Sunny tomorrow.
- 30% chance it will Rain tomorrow.
- 10% chance it will be Cloudy tomorrow.

Each row tells you the probabilities for the **next state**, given the **current state**.

For the above matrix, the stationary distribution  $\pi$  satisfies  $\pi P = \pi$   
After solving, you might find  $\pi = [0.5, 0.3, 0.2]$  which means in the long run,  
• 50% of the days will be Sunny.  
• 30% will be Rainy.  
• 20% will be Cloudy.

#### 8. Metropolis Hastings Algorithm

Let  $b(j)$ ,  $j = 1, \dots, m$  be positive numbers, and let  $B = \sum_{j=1}^m b(j)$ . Suppose that  $m$  is large and  $B$  is difficult to calculate, and that we want to simulate a sequence of random variables with probability mass function  $\pi(j) = \frac{b(j)}{B}$ . One way of simulating a sequence of random variables whose distributions converge to  $\pi(j)$  is to find a Markov chain that is easy to simulate and whose limiting probabilities are  $\pi(j)$ . Convergence is guaranteed as long as  $\pi(j) \propto b(j)$ .

- Start with an **irreducible proposal distribution**  $q(x, x')$ , which generates candidate samples  $X'$  from the current state  $X$
- $q(x, x')$  does not need to match the target distribution; it can be any valid probability function.
- A candidate sample  $X'$  is either **accepted** or **rejected** based on an **acceptance probability**:  $\alpha = \min(1, \frac{b(X')q(X, X')}{b(X)q(X', X)})$
- This ensures the Markov Chain converges to the target distribution  $\pi(x)$
- Generate a candidate sample  $X'$  using  $q(X, X')$
- Compute  $\alpha$  based on the target and proposal distributions.
- Draw a random number  $U \sim U[0, 1]$
- If  $U < \alpha$ , accept  $X'$  as the next state.
- Otherwise, stay at the current state  $X$
- Iterate this process to build a Markov Chain with  $\pi(x)$  as its stationary distribution.

#### 8.1 Choosing a proposal distribution

$q(x, x')$  is a **probabilistic function** that specifies how to propose a new candidate state  $X'$ , given the current state  $X$ .  $q(x, x')$  doesn't need to match the target distribution  $\pi(x)$ . However, it affects how efficient the algorithm is.

**Gaussian Proposal** (continuous space):  $q(x, x') = \frac{1}{\sqrt{2\sigma^2}} \exp\left(-\frac{(x' - x)^2}{2\sigma^2}\right)$ . Here,  $X'$  is drawn from a normal distribution centred at  $X$  with variance  $\sigma^2$ . Generate  $X' = X + \mathcal{N}(0, \sigma^2)$

**Uniform Proposal** (discrete space):  $q(x, x') = \begin{cases} \frac{1}{\text{neighbors of } x} & \text{if } x' \text{ is a valid neighbor of } x, \\ 0 & \text{otherwise.} \end{cases}$ . Generate,  $X'$  by randomly choosing a neighbour of  $X$ .

**Exploration:**  $q(x, x')$  must allow the Markov chain to explore the entire space  $S$ . Otherwise, the chain might "get stuck" in one part of the space.  
**Efficiency:**  $q(x, x')$  should balance between: Proposing states close to  $X$  (so  $\alpha$  is high) and exploring the state space widely.

#### Common Choices:

- Gaussian Proposal:** Works well for continuous variables.
- Uniform Proposal:** Simple for discrete spaces.
- Adaptive Proposal:** Adjust  $q(x, x')$  dynamically based on past samples to improve efficiency.

#### 9. Gibbs Sampler Algorithm

$$f(x|y, z) = \frac{f(x, y, z)}{f(y, z)} \propto f(x, y, z)$$

Conditional Marginal is Proportional to the Joint Posterior because the denominator is constant with respect to  $x$ . Therefore, we just need to take out the multiplicative constants (because of proportionality) to find the marginal conditional posterior densities. We use the most updated values to update the rest.

- Initialise the first guess  $\theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_d^{(0)})$ , where  $\theta$  represents the vector of all variables.
- For each iteration  $t$ , update each variable sequentially
- Sample  $\theta_1^{(t)}$  from its conditional distribution  $\pi(\theta_1 | \theta_2^{(t-1)}, \dots, \theta_d^{(t-1)}, Y)$
- Sample  $\theta_2^{(t)}$  from its conditional distribution  $\pi(\theta_2 | \theta_1^{(t)}, \theta_3^{(t-1)}, \dots, \theta_d^{(t-1)}, Y)$
- Continue this for all variables  $\theta_1, \dots, \theta_d$
- Repeat sequential sampling  $T$  times/until convergence. Output sample set  $\{\theta_1^{(t)}, \dots, \theta_d^{(t)}\}_{t=1}^T$  approx. target joint distribution

- Relies on choosing  $L$ , which may be non-trivial if dependencies are long-range.

#### 11.4 Confidence Intervals for Case 2: Regenerative Methods

Applicable for stationary Markov chains. Identify a "regeneration state"  $l$  that the Markov chain revisits. Between consecutive visits to  $l$ , the process "restarts," making segments between visits **independent**. By leveraging these independent segments, we can reduce dependence issues.

Define  $\sigma_k$  as the time of the  $k - th$  visit to state  $l$ . Define segment lengths  $D_k = \sigma_k - \sigma_{k-1}$  (time between visits). Define segment sums  $H_k = \sum_{i=\sigma_{k-1}+1}^{\sigma_k} h(W_i)$ , where  $h(W_i)$  is the function of interest. Use the segments to estimate:  $E(H_k) = \frac{E(H_k)}{E(D_k)}$  by replacing expectations with sample means.

- Converts a dependent process into independent segments, allowing for the use of i.i.d. techniques.
- Requires finding a suitable regeneration state, which may not always exist or be easy to identify.

Feature	Window Methods	Regenerative Methods
<b>Applicability</b>	General stationary sequences	Stationary Markov chains with regeneration
<b>Core Idea</b>	Summing autocovariance within a fixed window	Using independent regenerative segments
<b>Dependency Handling</b>	Approximate (truncates at $L$ )	Exact (segments are i.i.d.)
<b>Ease of Use</b>	Relatively easy; requires choosing $L$	Requires identifying regeneration states
<b>Computation</b>	Summing covariances; fast for short-range	Segment extraction; depends on regeneration

11.5 Confidence Intervals for Case 3: Hypothesis Testing	Procedure
	- Define $X_t = h(W_t)$
Batch Means	- Use additional functions
Discard problematic batches to ensure convergence	$h_1, h_2, h_3$ and plot their batch means
	- Discard initial batches if unstable
	- Divide into $b_1$ and $b_2$ batches
	- Calculate variances $S_1^2, S_2^2$
	- Test $\frac{S_1^2}{S_2^2} \sim F_{b_1-1, b_2-1}$
Independent Runs for Smoothing	- Perform multiple runs from the same initial state
Ensure smoothness and verify convergence to equilibrium	- Compute smoothed estimates $E(X_t)$
	- Perform runs with different initial states $X_0$
Widely Separated Initialisations	- Plot $\{X_t\}$ over time
Estimate burn-in time and assess the stability of the process	- Identify $T$ when graphs appear similar

#### 12. Bootstrap Sampling

- Original dataset:**  $X = \{2, 4, 6, 8, 10\}$
  - Bootstrap sample 1:** Randomly sample with replacement:  $X_1^* = \{4, 6, 6, 8, 2\}$
  - Bootstrap sample 2:** Randomly sample again:  $X_2^* = \{10, 10, 6, 8, 4\}$
  - Bootstrap sample 3:** Randomly sample again:  $X_3^* = \{6, 8, 8, 2, 4\}$
  - Generate  $N = 1000$  bootstrap samples accordingly
- Each bootstrap sample has the same size as the original dataset ( $n = 5$ ). Test statistics should be calculated for each sample before being aggregated back to be the statistics for the sampling distribution.

#### 12.1 Bootstrap Statistics

- Given an observed dataset  $x_1, x_2, \dots, x_n$ .
- Sample mean:  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- Generate  $N$  bootstrap samples by sampling **with replacement** from the observed dataset.
- Each bootstrap sample  $X^*$  will have  $n$  elements.
- For each bootstrap sample  $X^*$ , calculate the mean :  $\bar{x}^* = \frac{1}{n} \sum_{i=1}^n X_i^*$
- The bootstrap estimate of MSE is defined as:  $MSE(\hat{F}) = \frac{1}{N} \sum_{i=1}^N (\bar{x}_i^* - \bar{x})^2$

General: Mean the bootstrap statistics to approximate population statistics