

Useful Distributions:	
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	$Var(X) = p(1 - p)$
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	$X \sim \text{Bin}(n, p)$
PMF	$P(X = k) = \binom{n}{k} p^k (1 - p)^{n - k}$ for $k = 0, 1, \dots, n$
Expectation	$E(X) = np$
Variance	$Var(X) = np(1 - p)$
Properties	If $X_1, \dots, X_n$ are i.i.d. with distribution $\text{Bernoulli}(p)$ , then $X_1 + \dots + X_n \sim \text{Bin}(n, p)$
Distribution	Geometric
Description	Number of Bernoulli trials to obtain the first success
Notation	$X \sim \text{Geometric}(p)$
PMF	$P(X = k) = p(1 - p)^{k - 1}$ for $k = 1, 2, 3, \dots$
Expectation	$E(X) = \frac{1}{p}$
Variance	$Var(X) = \frac{1 - p}{p^2}$
Distribution	Poisson
Description	The 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	$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}$ $k = 0, 1, 2, \dots, \lambda > 0$
Expectation	$E(X) = \lambda$
Variance	$Var(X) = \lambda$
Properties	When n is large and p is small, np is moderate, $\text{Bin}(n, p) \rightarrow \text{Poisson}(np)$
Distribution	Negative Binomial
Description	Number of Bernoulli trials to obtain r successes
Notation	$X \sim \text{NB}(r, p)$
PMF	$P(X = k) = \binom{k - 1}{r - 1} p^r (1 - p)^{k - r}$ for $k = 1, 2, 3, \dots$
Expectation	$E(X) = \frac{r(1 - p)}{p}$
Variance	$Var(X) = \frac{pr}{(1 - p)^2}$
Distribution	Multinomial Distribution
Description	n - Number of trials k - Number of mutually exclusive events
Notation	$X \sim \text{NB}(r, p)$
PMF	$\frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$
Expectation	$E(X_i) = np_i$
Variance	$Var(X_i) = np_i(1 - p_i)$ $Cov(X_i, X_j) = -np_i p_j (i \neq j)$
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	$F(x) = \begin{cases} \frac{x - a}{b - a}, & a \leq x < b \\ 1, & b \leq x \end{cases}$

Expectation	$E(X) = \frac{a + b}{2}$
Variance	$Var(X) = \frac{(b - a)^2}{12}$
Properties	$U(0, 1) \equiv \text{Beta}(1, 1)$ <b>Transform to Uniform(a, b) from U(0,1):</b> $Y = (b - a)X + a$
Distribution	Normal
Description	$X \sim N(\mu, \sigma^2)$
PDF	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$ $-\infty < x < \infty$
CDF	$F(x) = \int_{-\infty}^x f(x) dx$ $-\infty < x < \infty$
Expectation	$E(X) = \mu$
Variance	$Var(X) = \sigma^2$
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 see the distribution of it and the normalising constant is straightforward
Distribution	Exponential
Description	$X \sim \text{Exp}(\lambda)$
PDF	$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$ Note that $\lambda > 0$
CDF	$F(x) = \begin{cases} 1 - e^{-\lambda x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$
Expectation	$E(X) = \frac{1}{\lambda}$
Variance	$Var(X) = \frac{1}{\lambda^2}$
Properties	For any $X \sim \text{Exp}(\lambda)$ $P(X > s + t   X > s) = P(X > t)$
Distribution	Cauchy
Description	$X \sim \text{Cauchy}(x_0, \gamma)$
Notation	$f(x) = \frac{1}{\pi\gamma} \left( \frac{1}{1 + \left(\frac{x - x_0}{\gamma}\right)^2} \right)$ $-\infty < x < \infty$
PDF	$F(x) = \frac{1}{\pi} \arctan\left(\frac{x - x_0}{\gamma}\right) + \frac{1}{2}$ $-\infty < x < \infty$
Inverse PDF	$F^{-1}(u) = \gamma \tan[\pi(u - 0.5)] + x_0$ $u \in [0, 1]$
Distribution	Standard Cauchy
Description	$X \sim \text{Cauchy}(0, 1)$
Notation	$f(x) = \frac{1}{\pi} \left( \frac{1}{1 + x^2} \right)$ $-\infty < x < \infty$
PDF	$F(x) = \frac{1}{\pi} \arctan x + \frac{1}{2}$ $-\infty < x < \infty$
Inverse PDF	$F^{-1}(u) = \tan[\pi(u - 0.5)]$ $u \in [0, 1]$
Distribution	Gamma
Notation	$X \sim \text{Gamma}(a, b)$
PDF	$g(x) = \begin{cases} \frac{\lambda^a}{\Gamma(\lambda)} x^{a - 1} e^{-\lambda x}, & x \geq 0 \\ 0, & t < 0 \end{cases}$
CDF	$G(x) = \frac{1}{\Gamma(\alpha)} \gamma(\alpha, \lambda x)$
Expectation	$E(X) = \frac{\alpha}{\lambda}$
Variance	$Var(X) = \frac{\alpha}{\lambda^2}$
Properties	<b>Gamma Function:</b> $\Gamma(z) = \int_0^\infty t^{z - 1} e^{-t} dt$ <b>Recursion Property:</b> $\Gamma(a + 1) = a\Gamma(a)$ <b>Gamma Function Computation:</b> $\Gamma(x) = (x - 1)!$ <b>Sum of Gamma Random Variables:</b> Let $X_1, \dots, X_k$ be independent random variables, assume $X_i \sim \text{Gamma}(a_i, b)$ for each $i$ . Then $X_1 + \dots + X_k \sim \text{Gamma}(a_1 + \dots + a_k, b)$

	<b>Connection with Standard Normal:</b> If $Z \sim N(0, 1)$ then $Z^2 \sim \text{Gamma}\left(\frac{1}{2}, \frac{1}{2}\right) \sim \chi^2(1)$ <b>Connection with Chi Squared:</b> Assume $Z_1, \dots, Z_k$ are i.i.d $N(0, 1)$ random variables. Then $Z_1^2 + \dots + Z_k^2 \sim \text{Gamma}\left(\frac{k}{2}, \frac{1}{2}\right) \sim \chi^2(k)$ <b>Connection with Exponential Distribution:</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 $X \sim \text{Gamma}(a, b)$ then $\lambda X \sim \text{Gamma}\left(a, \frac{b}{\lambda}\right)$
Distribution	Beta
Notation	$X \sim \text{Beta}(a, b)$
PDF	$f(x) = \frac{\Gamma(a + \beta)}{\Gamma(a)\Gamma(\beta)} x^{a - 1} (1 - x)^{b - 1}$ $0 \leq x \leq 1, \quad a > 0, b > 0$
CDF	$G(x) = \frac{1}{\Gamma(a)} \gamma(a, \lambda x)$
Expectation	$E(X) = \frac{a}{a + b}$
Variance	$Var(X) = \frac{ab}{(a + b)^2(a + b + 1)}$
Properties	<b>Swap of parameters:</b> If $X \sim \text{Beta}(a, b)$ , then $1 - X \sim \text{Beta}(b, a)$ If $X \sim \text{Gamma}(a, \beta)$ , $Y \sim \text{Gamma}(b, \beta)$ and $X, Y$ are independent, then $\frac{X}{X + Y} \sim \text{Beta}(a, b)$ <b>Order Statistics:</b> If $X_1, \dots, X_n$ are i.i.d from $\text{Uniform}(0, 1)$ and $X_{(1)} \leq \dots \leq X_{(n)}$ are their order statistics, then for $k = 1, \dots, n$ $X_{(k)} \sim \text{Beta}(k, n + 1 - k)$ <ul style="list-style-type: none"><li>Useful to know when we want to generate Beta distribution, we can just draw iid uniform and order them then pick the kth one</li></ul>

**Transformation of Random Variables**  
**Scaling and Shifting of Random Variables:** Suppose that  $X$  is a continuous random variable with pdf  $f(x)$ 

- Shift:** If  $a$  is a real number, then pdf of  $X + a$  is  $f(x - a)$
- Scale:** If  $b$  is a positive number, then the pdf of  $bX$  is  $b^{-1}f\left(\frac{x}{b}\right)$

**Change of variable formula:** Suppose  $U$  and  $V$  are functions of  $X$  and  $Y$ ,  $u = g_1(x, y)$ ,  $v = g_2(x, y)$ ,  $J(x, y) \neq 0$   
**Multivariable Joint Density of U and V:**  
 $f_{UV}(u, v) = f_{XY}(h_1(u, v), h_2(u, v))|J(x, y)|$   
Note that  $h_1(u, v)$  is  $x$  represented by  $u, v$  only.  $h_2(u, v)$  is  $y$  represented by  $u, v$  only.

**Jacobian:**  $J(x, y) = \det \begin{bmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{bmatrix} = \left( \frac{\partial u}{\partial x} \right) \left( \frac{\partial v}{\partial y} \right) - \left( \frac{\partial v}{\partial x} \right) \left( \frac{\partial u}{\partial y} \right)$   
Rows - Functions, Columns, Variables  
**Single Variable:** Suppose  $g(x)$  is a one-to-one differentiable function. If  $X$  has pdf  $f_X(x)$  and  $Y = g(X)$  then pdf of  $Y$  is:  
 $f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{dg^{-1}(y)}{dy} \right|$   
 $g^{-1}(y)$  is just  $x$  in terms of  $y$  and we substitute it into  $f_X$   
**Note** that if  $g(x, y)$  is not one-to-one, we break it into intervals such that it is one to one and we just add up the distribution on the range where they are one-to-one

**Sample Estimators:**  
Let  $X_1, X_2, \dots$  be a sequence of iid random variables with mean  $\mu$  and variance  $\sigma^2$   
**Sample Mean:** The mean of the sample that we are currently looking at  
 $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i, \quad E(\bar{X}_n) = \mu$

**Sample Variance:** Variance of the sample data that we are currently looking at  
 $S_n^2 = \frac{1}{n - 1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad E(S_n^2) = \sigma^2$   
**Variance of Sample Mean:** Variation of the sample means that we will get over the  $n$  samples  
 $Var(\bar{X}_n) = \frac{\sigma^2}{n}$

**Central Limit Theorem (CLT iid version):** Suppose that the random variable  $X$  has finite second moment (i.e.  $E[X^2] < \infty$ ), then the following **convergence in distribution** holds  
 $\lim_{n \rightarrow \infty} \sqrt{n}(\bar{X}_n - \mu) = N(0, \sigma^2)$ 

- $\bar{X}_n - \mu$  converges to 0 in order of  $n^{-\frac{1}{2}}$
- In the multivariate case, replace  $\sigma^2$  by the covariance matrix

**Covariance:**  
 $Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$   
 $= E(XY) - E(X)E(Y)$   
 $Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$

**Discrete Random Variable Generation:**  
Compute the probability for each of the possible values for the pmf. Then we just generate a uniform distribution to check which of the probability range it lies within.  
**Algorithm:**

- Generate  $U \sim \text{Uniform}(0, 1)$
- If  $U < p_0$ , set  $X = x_0$  and stop
- If  $U < p_0 + p_1$ , set  $X = x_1$  and stop
- :
- Otherwise, set  $X = x_n$

**Inversion Method (Continuous Random Variable):**  
For a given random variable  $X$ , if we want to generate it, we can do the following. If we are given the pdf of the random variable  $X$ ,  $f(x)$ :

- Integrate  $f(x)$  over the entire range to get the CDF,  $F(x)$
- Let  $U \sim \text{Uniform}(0, 1)$ . Set  $U = F(x)$  and find the inverse of the cdf  $F^{-1}(u) = X$
- Once we have found the inverse CDF, we can just generate a uniform distribution and put inside the inverse CDF to get one  $X$

**Algorithm:**

- Generate  $U \sim \text{Uniform}(0, 1)$
- Set  $X = F^{-1}(U)$

**Note:**

- For  $\text{Exp}(\lambda)$ , we can use the following to generate the random variable:  
 $X = -\frac{1}{\lambda} \log U$
- For  $\text{Gamma}(n, \lambda)$ , we can use the following to generate the random variable:  
Note that we are making use of the fact that  $\text{Gamma}(n, \lambda)$  is the sum of  $n$   $\text{Exp}(\lambda)$   
 $X = -\frac{1}{\lambda} \log(U_1 \dots U_n)$

**Rejection Sampling:**  
**Quick ways to check that  $\sup_x \frac{f(x)}{g(x)} < +\infty$**  (But still need to rigorously show, this can give a brief idea)

- Domain of  $g(x)$  should cover the domain of  $f(x)$
- The tails of the proposal  $g(x)$  should be heavier than the tails of  $f(x)$

**Rigorous ways to check:**

- Differentiate the ratio of  $\frac{f(x)}{g(x)}$  and find the maximum value that the ratio can attain
- Try to observe what will happen to the ratio  $\frac{f(x)}{g(x)}$  when  $x \rightarrow +\infty$ . Look at what kind of function it will look like and make the conclusion from there

Theoretical number of simulations required to get 1 acceptance:  $M = \sup_x \frac{f(x)}{g(x)}$   
**Logical steps to do (When we are computing):**

- Try to imagine the shape of  $g(x)$  and  $f(x)$ , when one increases, the other should increase also, vice versa
- Find the value of  $M = \sup_x \frac{f(x)}{g(x)}$ . Check that it exists and state the value where we can compute the maximum value using the rigorous way to check
- Specify the rejection function of  $\frac{f(x)}{Mg(x)}$  and our  $U$  needs to be within the rejection function range else it will be rejected
- Generate  $Y$  using some kind of method (normally inversion)

**Algorithm:**

- Generate  $Y \sim g$
- Generate  $U \sim \text{Uniform}(0, 1)$
- If  $U \leq \frac{f(Y)}{Mg(Y)}$ , then accept: set  $X = Y$  and stop. Otherwise, reject and return to step 1

**Unknown Normalising Constant:**  
If we only know  $f(x)$  up till a certain normalising constant, it will work the same, just take the ratio and supremum to be:  
 $\frac{\tilde{f}(x)}{g(x)}, \quad \sup_x \frac{\tilde{f}(x)}{Mg(x)}$

**Polar Method for Bivariate Normal:**  
 $S = R^2 = X^2 + Y^2, \tan \theta = \frac{Y}{X}, X = R \cos(\theta), Y = R \sin(\theta)$   
**Change of variable from (X, Y) to (S, \theta)**  
 $f(s, \theta) = \frac{1}{2} e^{-\frac{s}{2}} \frac{1}{2\pi}, \quad 0 < s < \infty, 0 < \theta < 2\pi$   
 $S = R^2 \sim \text{Exp}\left(\frac{1}{2}\right)$  and  $\theta \sim \text{Uniform}(0, 2\pi)$

**Box-Muller Algorithm v1:**

- Generate random numbers  $U_1 \sim \text{Uniform}(0, 1)$  and  $U_2 \sim \text{Uniform}(0, 1)$
- Set:  
 $X = \sqrt{-2 \log U_1} \cos(2\pi U_2)$   
 $Y = \sqrt{-2 \log U_1} \sin(2\pi U_2)$

**Box-Muller Algorithm v2:** Suppose that  $(V_1, V_2)$  is uniformly distributed in the disk centered at  $(0, 0)$  with radius 1 and the random angle is  $\theta \sim \text{Uniform}(0, 2\pi)$ 

- Generate random numbers  $U_1 \sim \text{Uniform}(0, 1)$  and  $U_2 \sim \text{Uniform}(0, 1)$
- Set  $V_1 = 2U_1 - 1, V_2 = 2U_2 - 1, S = V_1^2 + V_2^2$  ( $V_1, V_2$  are just  $X$  and  $Y$  coordinates sampled from  $\text{Uniform}(-1, 1)$ )
- If  $S > 1$ , return to Step 1 ( $S$  is the radius squared for a unit disk so it should be  $\leq 1$ )
- Return the independent unit normals  
 $X = \sqrt{\frac{-2 \log S}{S}} V_1, \quad Y = \sqrt{\frac{-2 \log S}{S}} V_2$

**Simple Sampling:** Sample  $X_1, X_2, \dots, X_n$  independently from  $f$ , we can estimate the true parameter shown below by  
 $\theta = E[\varphi(X)] = \int_S \varphi(x) f(x) dx$

**Simple Sampling Estimator:**  $\hat{\theta}_{SS} = \frac{1}{n} \sum_{i=1}^n \varphi(X_i)$   
**Simple Sampling Exact Variance of  $\hat{\theta}$  (Variance of Sample Mean):**  $Var(\hat{\theta}) = \frac{Var(\varphi(X))}{n} = \frac{\int_S \varphi^2(x) f(x) dx - \theta^2}{n}$   
**Simple Sampling Asymptotic Variance of  $\hat{\theta}$ :**  
 $\sigma^2 \equiv Var[\varphi(X)] = \int_S \varphi^2(x) f(x) dx - \theta^2$

**Simple Sampling Estimated Asymptotic Variance of  $\hat{\theta}$ :**  
**Note that this is not an unbiased estimate of  $\sigma^2$**   
 $\hat{\sigma}_{SS}^2 = \frac{1}{n} \sum_{i=1}^n \varphi^2(X_i) - \hat{\theta}_{SS}^2$

**Simple Sampling Estimated Variance of  $\hat{\theta}$  (Sample Variance):**  $\widehat{Var}(\hat{\theta}) = \frac{\hat{\sigma}_{SS}^2}{n}$

## Simple Sampling Asymptotic Confidence Interval for $\theta$ :

$$\left[ \hat{\theta} - 1.96 \frac{\hat{\sigma}_{\text{SS}}}{\sqrt{n}}, \hat{\theta} + 1.96 \frac{\hat{\sigma}_{\text{SS}}}{\sqrt{n}} \right]$$

## Importance Sampling:

Sample  $X_1, X_2, \dots, X_n$  independently from  $g$ , we can estimate the true parameter shown below by

$$\hat{\theta} = E_g[\varphi(X)] = \int_S \varphi(x) f(x) dx$$

$$\hat{\theta} = \int_S \frac{\varphi(x) f(x)}{g(x)} g(x) dx = E_g \left[ \frac{\varphi(x) f(x)}{g(x)} \right] = E_g[\varphi(Y) w(Y)]$$

**Weighting Function:**  $w(y) = \frac{f(y)}{g(y)}$

**Importance Sampling Estimator:** Unbiased estimator of  $\theta$

$$\hat{\theta}_{\text{IS}} = \frac{1}{n} \sum_{i=1}^n \frac{\varphi(x_i) f(x_i)}{g(x_i)} = \frac{1}{n} \sum_{i=1}^n \varphi(x_i) w(x_i)$$

**Importance Sampling Exact Variance of  $\hat{\theta}$  (Variance of Sample Mean):**

$$\text{Var}(\hat{\theta}) = \frac{\text{Var}[\varphi(X) w(X)]}{n} = \frac{\int_S \frac{\varphi^2(x) f^2(x)}{g(x)} dx - \theta^2}{n}$$

**Importance Sampling Asymptotic Variance of  $\hat{\theta}$ :**

$$\sigma_{\text{IS}}^2 \equiv \text{Var}[\varphi(X) w(X)] = \int_S \frac{\varphi^2(x) f^2(x)}{g(x)} dx - \theta^2$$

**Importance Sampling Estimated Asymptotic Variance of  $\hat{\theta}$ :** Note that this is not an unbiased estimate of  $\sigma^2$

$$\hat{\sigma}_{\text{IS}}^2 = \frac{1}{n} \sum_{i=1}^n \frac{\varphi^2(x_i) f^2(x_i)}{g(x_i)} - \hat{\theta}_{\text{IS}}^2$$

**Importance Sampling Estimated Variance of  $\hat{\theta}$  (Sample Variance):**  $\hat{\text{Var}}(\hat{\theta}) = \frac{\hat{\sigma}_{\text{IS}}^2}{n}$

**Importance Sampling Asymptotic Confidence Interval for  $\theta$ :**  $\left[ \hat{\theta} - 1.96 \frac{\hat{\sigma}_{\text{IS}}}{\sqrt{n}}, \hat{\theta} + 1.96 \frac{\hat{\sigma}_{\text{IS}}}{\sqrt{n}} \right]$

**Optimal  $g$ :**  $g(x) \propto |\varphi(x)| \cdot f(x)$

- Asymptotic variance of  $\hat{\theta}_{\text{IS}}$  with the proposal density  $g$  is exactly 0 if  $\varphi(x) \geq 0$  for all  $X \in S$

**To find  $g(x)$ :**

- Let  $h(x) = c|\varphi(x)|f(x)$
- Let  $1 = \int_S h(x) = \int_S c|\varphi(x)|f(x)$  and solve for  $c$  (Note that if we have the value for  $I = \int_S |\varphi(x)|f(x) dx \Rightarrow c = \frac{1}{I}$ )
- $g(x) = ch(x)$

**Self-Normalizing Importance Sampling:** We only know the distribution of  $f$  and  $g$  up to a normalising constant ( $Z_f > 0, Z_g > 0$ )

$$\tilde{f}(x) = \frac{f(x)}{Z_f}, \quad \tilde{g}(x) = \frac{g(x)}{Z_g}$$

**Generalised weights:**  $\tilde{w}(x) = \frac{\tilde{f}(x)}{\tilde{g}(x)}$ , for all  $x \in S$

**Self-normalised importance sampling estimator of  $\theta$ :**  $E_f[\varphi(X)] = \int_S \varphi(x) f(x) dx$

$$\hat{\theta}_{\text{SIS}} = \frac{\sum_{i=1}^n \varphi(X_i) \tilde{w}(X_i)}{\sum_{i=1}^n \tilde{w}(X_i)}$$

**Asymptotic variance of  $\hat{\theta}_{\text{SIS}}$ :** Normally larger than the IS version because of the random denominator

$$\hat{\sigma}_{\text{SIS}}^2 = E_g(w^2(X) \cdot [\varphi(X) - \theta]^2)$$

Where  $w(x) = f(x)/g(x)$ , it is the true weight

**Exact Variance for  $\hat{\theta}_{\text{SIS}}$ :** No closed form

- Note that  $\text{Var}[\hat{\theta}_{\text{SIS}}] \neq \frac{\hat{\sigma}_{\text{SIS}}^2}{n}$  (Not the same as simple sampling and importance sampling)

**Estimator of the Variance of  $\hat{\theta}_{\text{SIS}}$**

$$\hat{\text{Var}}(\hat{\theta}_{\text{SIS}}) = \frac{\hat{\sigma}_{\text{SIS}}^2}{n} = \frac{\sum_{i=1}^n [\tilde{w}^2(X_i) \varphi(X_i) - \hat{\theta}_{\text{SIS}}]^2}{\sum_{i=1}^n \tilde{w}^2(X_i)}$$

**95% Asymptotic Confidence Interval**

$$\left[ \hat{\theta} - 1.96 \sqrt{\frac{\hat{\sigma}_{\text{SIS}}^2}{n}}, \hat{\theta} + 1.96 \sqrt{\frac{\hat{\sigma}_{\text{SIS}}^2}{n}} \right]$$

**Calculus Results:**

- $\int_1^{+\infty} \frac{1}{x^p} = \frac{1}{-p+1} x^{-p+1} \Big|_1^{+\infty} = \begin{cases} < +\infty & \text{if } p > 1 \\ +\infty & \text{if } p \leq 1 \end{cases}$

- $\int_0^1 \frac{1}{x^p} = \frac{1}{-p+1} x^{-p+1} \Big|_0^1 = \begin{cases} < +\infty & \text{if } p < 1 \\ +\infty & \text{if } p \geq 1 \end{cases}$
- $\int_0^1 e^x = \begin{cases} < +\infty & \text{if } p < 1 \\ +\infty & \text{if } p \geq 1 \end{cases}$

**Rare Event Estimation:** When the  $p^*$  we want to estimate is small

**Relative Standard Deviation** =  $\frac{\text{asymptotic s.d.}}{p^*}$

- Checks the magnitude of the asymptotic sd of our estimator as compared to the actual value  $\rightarrow$  If it is large, it means that the magnitude of the sd of our estimator is larger than the actual value and it will give a very bad estimate

For a Bernoulli RV  $\rightarrow$  **Relative s.d.** =  $\frac{\sqrt{p(1-p)}}{p} = \sqrt{\frac{1-p}{p}}$

(Therefore, if the probability of it happening is low then the sd is high)

- Consider using a density centered at the point where we need more points so that the probability is higher and lowering the relative sd
- Remember that we can take the  $\varphi$  as the indicator function to indicate  $P(X_i > 4)$  for example.

**Control Variates Method:** Using  $\hat{h}$  correlated with  $\hat{\theta}$

$$\hat{\theta} = \hat{\theta} + \beta[\hat{h} - E_f[h(x)]]$$

$$\text{Var}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \beta^2 \text{Var}(\hat{h}) + 2\beta \text{Cov}(\hat{\theta}, \hat{h})$$

$$\hat{\theta}^* = -\frac{\text{Cov}(\hat{\theta}, \hat{h})}{\text{Var}(\hat{h})}, \quad \text{Var}^*(\hat{\theta}) = (1 - \text{Cor}(\hat{\theta}, \hat{h})^2) \text{Var}(\hat{\theta})$$

$$\text{Cor}(\hat{\theta}, \hat{h}) = \text{Cov}(\hat{\theta}, \hat{h}) / \sqrt{\text{Var}(\hat{\theta}) \text{Var}(\hat{h})}$$

**Antithetic Variates Method**

**Useful Facts:**

- If  $X$  is generated from an inversion method from a cdf  $F$ , i.e.  $X = F^{-1}(U)$  with  $U \sim \text{Uniform}(0, 1)$ , then  $X' = F^{-1}(1 - U)$  also follows distribution  $F$ . Note that  $F^{-1}$  is a monotonic function
- If we want to compare any other estimator with the Antithetic Estimator, remember to take sample size as  $2n$
- For  **$U(0, 1)$**   **$U$**  and  **$1 - U$**  are perfectly negatively correlated. For  **$N(0, 1)$** ,  **$X$**  and  **$-X$**  are perfectly negatively correlated. We just need to find a monotonic function, either the inverse CDF or the function already given to us like  $e^x$ . Then it will be  $h(X)$  and  $h(1 - X)$  or whatever that is perfectly negatively correlated.

**Antithetic Estimator:**

$$\hat{I}_{\text{an}} = \frac{1}{2n} \sum_{i=1}^n (h(U_i) + h(1 - U_i))$$

**Variance of Antithetic Estimator:**

$$\text{Var}(\hat{I}_{\text{an}}) = \frac{1}{2n} [\text{Var}(h(U)) + \text{Cov}(h(U), h(1 - U))]$$

**Expectation-Minimization (EM) Algorithm:**

**Jensen Inequality:**  $f(E[X]) \leq E[f(X)]$

**Important Result:** For any positive function  $\varphi(x)$  and a density  $q(x)$

$$\log \left( \int \varphi(x) q(x) dx \right) \geq \int \log[\varphi(x)] q(x) dx$$

**Likelihood Function:**  $L(Y|\theta)$  where  $Y$  is the observed data and  $\theta$  are the parameters

**Log-likelihood Function:**  $\ell(Y|\theta) = \log L(Y|\theta)$

**Complete data log-likelihood:**

$$\ell^c(Y, Z|\theta) = \log p(Y, Z|\theta) = \sum_{i=1}^n \log p(y_i, z_i|\theta)$$

**Steps:**

- Expectation:** Given  $\theta_k$ , calculate the function  $Q(\theta|\theta_k) = E_{z \in Z}[\ell^c(Y, Z|\theta)|Y, \theta_k]$
- Maximisation:** Calculate the next iterate of the parameters  $\theta$

$$\theta_{k+1} = \arg \max_{\theta \in \Theta} Q(\theta|\theta_k)$$

- Differentiate the  $Q(\theta|\theta_k)$  function to get the maximiser

**Algorithm:**

- Initialize  $\theta_0 = (\mu_1^{(0)}, \mu_2^{(0)}, \dots)$  and  $\epsilon > 0$
- E-Step:** In the  $k$ th iteration (given  $\theta_k = (\mu_1^{(k)}, \mu_2^{(k)}, \dots)$ ) calculate  $q_i^{(k,j)}$  for  $i = 1, \dots, n$ ,  $j = 1, \dots, m$ . Where  $i$  is the  $i$ th observation and  $j$  is the  $j$ th parameter
- M-Step:** Calculate  $\theta_{k+1} = (\mu_1^{(k+1)}, \mu_2^{(k+1)}, \dots)$  using (2)
- Iterate between the E-step and M-step until convergence.  $\{|\theta_{k+1} - \theta_k| < \epsilon\}$

**Tips:**

- For mixture distribution where we do not know the relative number of samples in each group, we can make use of an indicator  $I_{\{z_i=1\}}$  to indicate which density to use.

Note that once we look at the  $Q$  function, this will be  $E_z[I\{z_i = 1\}Y, \theta_k]$  and it becomes

$$\begin{aligned} q_i^{(k,1)} &= p(z_i = 1|y_i, \theta_k) \\ &= \frac{p^{(k)}(y_i|z_i = 1, \theta_k)}{p^{(k)}(y_i|z_i = 1, \theta_k) + (1 - p^{(k)})p(y_i|z_i = 2, \theta_k)} \\ q_i^{(k,2)} &= 1 - q_i^{(k,1)} \end{aligned}$$

If there is more than 1 latent variable, then it will still be the same idea. Just that we cannot take  $1 - q_i^{(k,2)}$  and denominator, we will sum over all possible combinations of  $z_i$  using the idea of Law of Total Probability

- If we know the number of relative samples in each group, like the truncated Poisson, we may just need to take expectation over those  $z_i$  terms and for truncated Poisson it will be a **negative binomial distribution**. (Can look at the mean of it when we are computing expectation of  $E[z|Y, \lambda_k]$ )

**Markov Chains**

**Stationary Distribution:**  $\pi P = \pi$  where  $\pi$  is the stationary distribution and  $P$  is the transition matrix

- When finding the stationary distribution, we can consider the  $\pi$  of those states that are recurrent. Those that are transient will have probability 0 at the stationary distribution.

**Transition Matrix  $P$ :** Tells us the probability of going from a location to the next. Every row is the starting position and the columns are the position that we are going next.

**Nice Properties:**

- For nice Markov Chains (Irreducible, Positive Recurrent, Aperiodic), in the long run, it will converge to its stationary distribution.
- The empirical distribution of  $\{X_1, X_2, \dots, X_T\}$  as  $T \rightarrow \infty$  (Noting that they are dependent samples since it is a Markov Chain) is close to the stationary distribution.

**Properties of Stationary Distributions:**

- The stationary distribution  $\pi$  exists and is unique, if the Markov chain is irreducible and positive recurrent
- $\lim_{t \rightarrow \infty} P^t = 1\pi$  where  $1 = (1, \dots, 1)^T$  holds if the Markov chain is irreducible, positive recurrent and aperiodic
- To find the stationary distribution  $\pi$ , solve the system of equation  $\pi P = \pi$  or use the detailed balance condition.
- $\lim_{t \rightarrow \infty} P_{ab}(t) = \pi_b$  where  $\pi$  is the stationary distribution
- If for each row, the probability is equal for any non-zero entries. Then the stationary distribution for  $\pi_i$  will be the number of neighbours it has / total number of neighbours for everyone. (i.e.  $\frac{1}{2}$  for row 1 and  $\frac{1}{3}$  for row 2,  $\pi_i = 2/5$ )

**Terminologies:**

- Irreducible Sets:** All states within this set are accessible between one another
- Positive Recurrent:** This happens if  $P[t_{ii} < \infty] = 1$  for some (and hence all) states  $i$ . This means that we can always come back to this state after a finite number of steps. (**Always for finite state space**)

- Transient:** If we are not able to come back to a state after a finite number of ways after we leave
- Aperiodic:** We should be able to come back to the same point in irregular moves

**Ergodic Theorem:** If  $X$  is irreducible and positive recurrent (**Only concerned with recurrent states**), and a function  $h(x)$  satisfies  $E_{\pi}[|h(x)|] < \infty$ , then

$$\frac{1}{N} \sum_{i=1}^N h(X_i) \rightarrow E_{\pi}[h(X)], \quad \text{as } N \rightarrow \infty, \text{ with probability } 1$$

Where  $E_{\pi}[h(X)] = \sum_i h(i)\pi_i$ , expectation of  $h(x)$  with respect to  $\pi(\cdot)$

**Bayesian Statistics:**

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior} \\ \pi(\theta|Y) \propto P(Y|\theta)\pi(\theta)$$

Note that the denominator for the usual Bayes Theorem is a constant in terms of  $\theta$  since we integrate out  $\theta$  so under the proportionality, we can remove it.

To find the posterior, just find the likelihood first and then multiply with the prior. Under proportionality, we can do that and try to find the form of it

**Useful Trick for Splitting:**

$$\begin{aligned} \sum_{i=1}^n (y_i - \theta)^2 &= \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - \theta)^2 \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\bar{y} - \theta)^2 + 2 \sum_{i=1}^n (y_i - \bar{y})(\bar{y} - \theta) \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \theta)^2 \end{aligned}$$

- $2 \sum_{i=1}^n (y_i - \bar{y})(\bar{y} - \theta) = 0$  because  $\sum_{i=1}^n (y_i - \bar{y}) = 0$  but not the case for the squared version
- $(\bar{y} - \theta)^2$  independent of  $i$ , therefore, we can just sum over  $n$  of them.

**Metropolis Hastings Algorithm:**

**Symmetric Transition Kernels:**

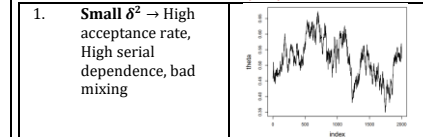
- Uniform kernel:**  $Q(\theta_{a-1})$  is the density of Uniform  $[\theta_a - \delta, \theta_a + \delta]$ , for some user-set  $\delta > 0$ .  
 $Q(\theta_a, \theta_b) = \frac{1}{2\delta}$ , for  $|\theta_b - \theta_a| < \delta$
- Normal kernel:**  $Q(\theta_{a-1})$  is the density of  $N(\theta_a, \delta^2)$ , for some user-set  $\delta^2 > 0$ .  
 $Q(\theta_a, \theta_b) = \frac{1}{\sqrt{2\pi}\delta} \exp\left\{-\frac{(\theta_b - \theta_a)^2}{2\delta^2}\right\}$

**Algorithm:**

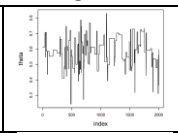
- Set  $\theta^{(0)}$  to some initial value  $\theta^{(0)} \in \Theta$
- For  $t = 1, \dots, T$ 
  - Draw  $\theta_t^* \sim N(\theta_{t-1}^{(t-1)}, \delta_t^2)$ ,  $\dots, \theta_n^*$  from  $N(\theta_{t-1}^{(t-1)}, \delta_t^2)$ . Could be common  $\delta^2$ . Check transition kernel formula
  - Compute the acceptance probability  
 $\alpha(\theta^{(t)}, \theta^*) = \min\left(1, \frac{p(Y|\theta^*)\pi(\theta^*)Q(\theta^*, \theta^{(t)})}{p(Y|\theta^{(t)})\pi(\theta^{(t)})Q(\theta^{(t)}, \theta^*)}\right)$
- Note that if the transition kernel is symmetric, then it will reduce to  
 $\alpha(\theta^{(t)}, \theta^*) = \min\left(1, \frac{p(Y|\theta^*)\pi(\theta^*)}{p(Y|\theta^{(t)})\pi(\theta^{(t)})}\right)$
- Generate  $U \sim \text{Uniform}(0, 1)$
- If  $U < \alpha(\theta^{(t-1)}, \theta^*)$ , then accept and set  $\theta^{(t)} = \theta^*$
- Otherwise reject  $\theta^*$  set  $\theta^{(t)} = \theta^{(t-1)}$
- Repeat until  $t = T$ . Return  $\{\theta^{(t)}\}_{t=1}^T$

**Trace Plots for Corresponding  $\delta^2$ :**

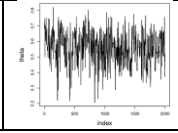
*Size is relative to the actual value of  $\delta^2$  for the data*



- Large  $\delta^2 \rightarrow$  Low acceptance rate, High serial dependence, bad mixing**



- Proper choice of  $\delta^2 \rightarrow$  Moderate acceptance rate, good mixing**



**Optimal Acceptance Rate:**

- For 3D and above: 0.234
- For 1D and 2D:  $\sim 30\%-50\%$

**Gibbs Sampler:**

**Computing the Marginal Conditional Posteriors**

$$\begin{aligned} f(x|y, z) &= \frac{f(x, y, z)}{f(y, z)} \\ &\propto f(x, y, z) \end{aligned}$$

- Conditional Marginal is Proportional to the Joint Posterior because the denominator is constant with rest 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.

**Algorithm:**

- Initialise  $\theta^{(0)} \in \Theta$  where  $\theta = (\theta_1, \dots, \theta_d)$
- At the step  $t$  ( $t = 1, \dots, T$ ), do the following sequentially
  - Sample  $\theta_1^{(t)} \sim \pi(\theta_1 | \theta_2^{(t-1)}, \theta_3^{(t-1)}, \dots, \theta_d^{(t-1)}, Y)$
  - Sample  $\theta_2^{(t)} \sim \pi(\theta_2 | \theta_1^{(t)}, \theta_3^{(t-1)}, \dots, \theta_d^{(t-1)}, Y)$
  - ...
  - Sample  $\theta_d^{(t)} \sim \pi(\theta_d | \theta_1^{(t)}, \theta_2^{(t)}, \dots, \theta_{d-1}^{(t)}, Y)$
- Repeat the steps above until time  $t = T$ . Output  $\{\{\theta_1^{(t)}, \dots, \theta_d^{(t)}\}_{t=1}^T\}$

**Dirichlet Distribution:**

$$f(x_1, x_2, \dots, x_d) \propto x_1^{a_1-1} x_2^{a_2-1} \dots x_d^{a_d-1}$$

- $0 < x_1 < 1, \dots, 0 < x_d < 1, \quad x_1 + x_2 + \dots + x_d = 1$

Degenerate distribution because only  $d - 1$  free points because  $x_d = 1 - x_1 - x_2 - \dots - x_{d-1}$

**Property:** They are marginally Beta Distributions

If  $(X_1, \dots, X_d) \sim \text{Dirichlet}(a_1, \dots, a_d)$

$$f(x_1, \dots, x_d) \propto x^{a_1-1} \dots x^{a_d-1}$$

Then for every  $i = 1, \dots, d$ ,  $X_i \sim \text{Beta}(a_i, \sum_{j \neq i} a_j)$

**Trick for Beta and Dirichlet Distribution:**

If we have multiple variables and we are unable to get the  $\theta(1 - \theta)$  form

$$f(x|y, z) \propto x^{a_1-1} (1 - x - y - z)^{a_d-1}$$

If we see this, just need to divide with  $1 - y - z$  (Everything other than the variable we are concerned about)

$$\begin{aligned} f(x|y, z) &\propto \left(\frac{x}{1 - y - z}\right)^{a_1-1} \left(\frac{1 - x - y - z}{1 - y - z}\right)^{a_d-1} \\ &= \frac{x^{a_1-1}}{1 - y - z} \sim \text{Beta}(a_1 + 1, a_d + 1) \\ &\propto (1 - Y - Z)U \end{aligned}$$

- Note that because this is a linear transformation, we can just multiply  $(1 - Y - Z)$  to  $U$  or else we need to do the change-of-variable formula

**Trick for Binomial**

- Try to find the  $\frac{$