

# Lecture 7: Monte Carlo Methods

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

April 30, 2024

# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

# Motivation I

If you can not calculate a probability or expectation exactly, then you have three powerful strategies:

- Simulations using Monte Carlo Methods
- Approximations using limiting theorems
  - ▶ Poisson approximation: The Law of Small Numbers
  - ▶ Sample mean limit: The Law of Large Numbers
  - ▶ Normal approximation: The Central Limit Theorem
- Bounds (upper and lower bounds) on probability using inequalities.

# Motivation II

Probability  
Math



Statistics  
Science

Monte Carlo  
Computing

# Monte Carlo Methods

- One of the top ten algorithms for science and engineering in 20th century
- Monte Carlo Methods, Simplex Method, Fast Fourier Transform, Quicksort, QR Algorithm...

# Widely Applications

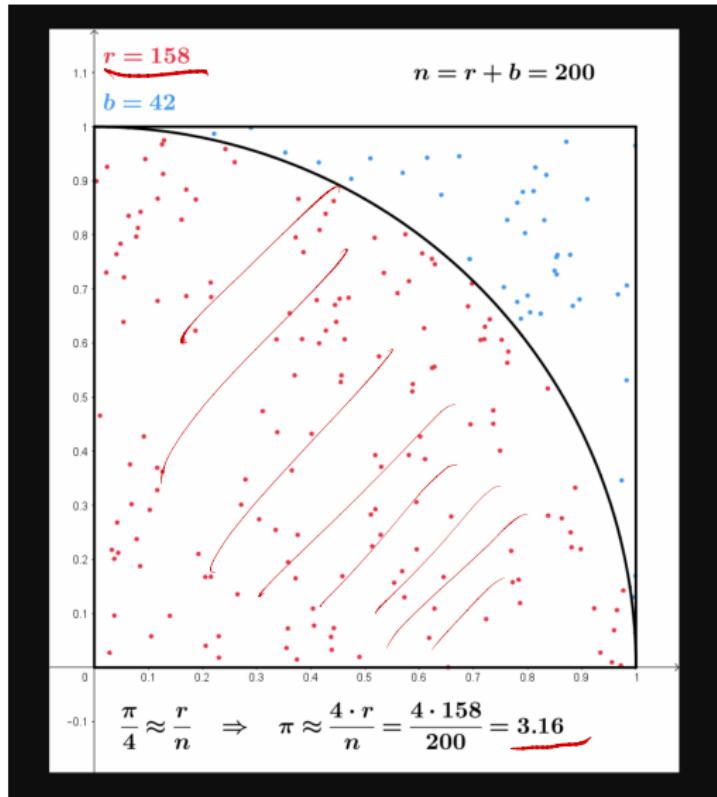
Monte Carlo methods have been used in various tasks, including

- Sampling from the underlying probability distribution  $f(x)$  and simulating a random system
- Sampling from posterior distribution for bayesian inference
- Estimation through numerical integration

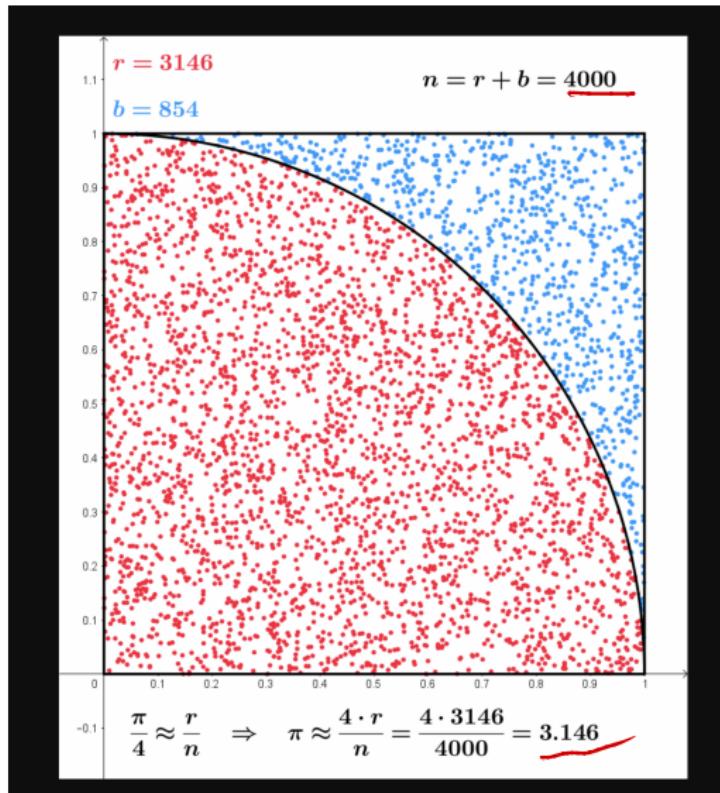
$$c = E_{\pi}(h(x)) = \int f(x)h(x)dx.$$

- Optimizing a target function to find its maxima or minima

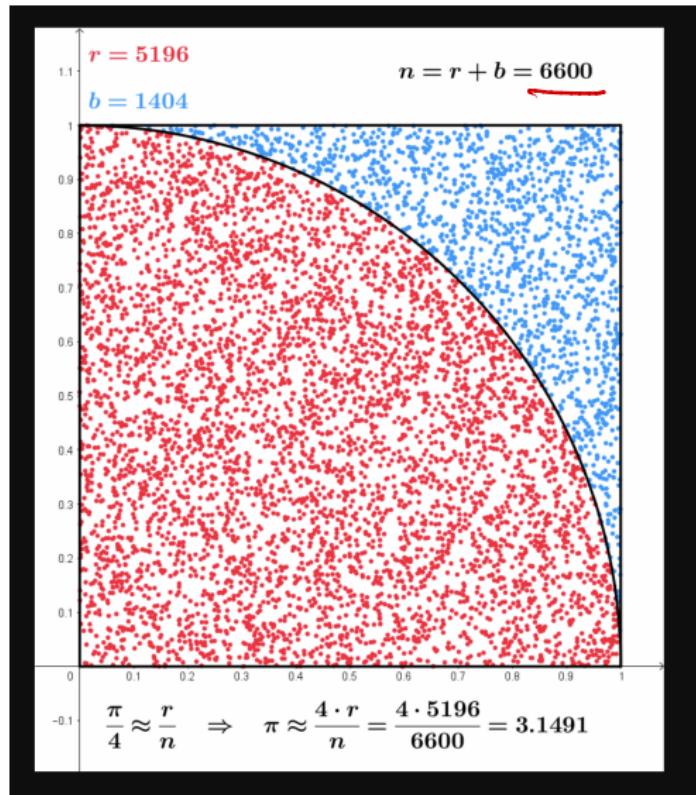
# Classical Example: Estimation of $\pi$



# Classical Example: Estimation of $\pi$



# Classical Example: Estimation of $\pi$



# History



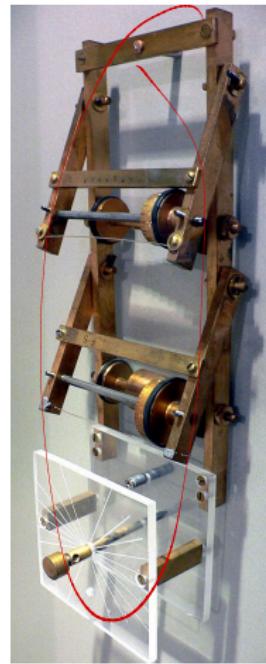
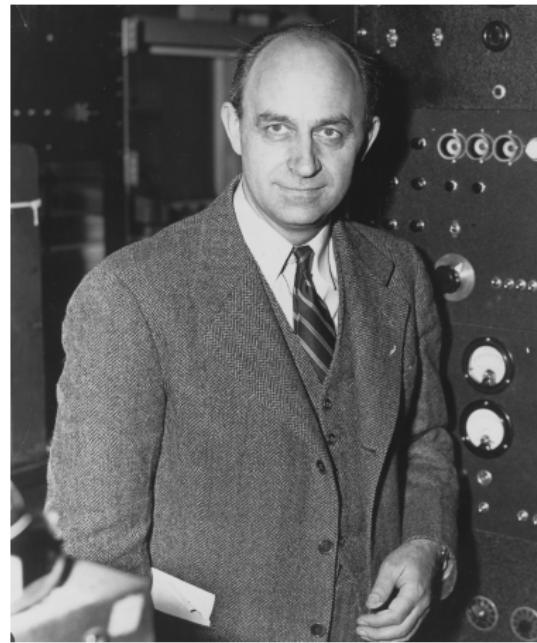
# Monte Carlo Methods

- Basic Monte Carlo methods: formally proposed by Stanislaw Ulam & John Von Neumann in 1940s at Los Alamos National Lab (Named after a casino in Monaco)



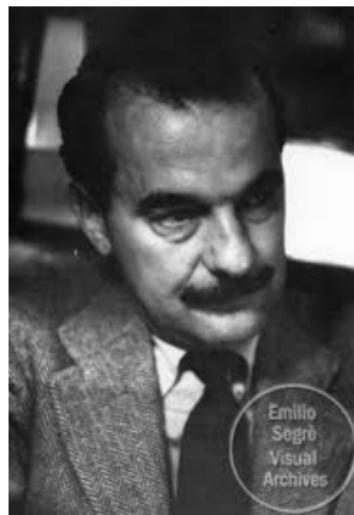
# Monte Carlo Trolley

- Analog computer invented by Enrico Fermi in 1946



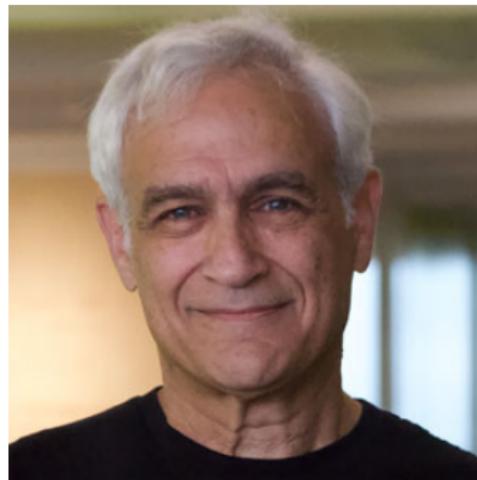
# Markov Chain Monte Carlo Methods

- Metropolis-Hastings Algorithm: formally proposed by Nicholas Metropolis et al in 1950s at Los Alamos National Lab, then extended in 1970 by Wilfred Keith Hastings



# Markov Chain Monte Carlo Methods

- Gibbs Sampling Algorithm: proposed in 1984 by brothers Stuart Geman (1949-) and Donald Geman (1943-).
- Gibbs sampling is named after the physicist Josiah Willard Gibbs (1839-1903), in reference to an analogy between the sampling algorithm and statistical physics.

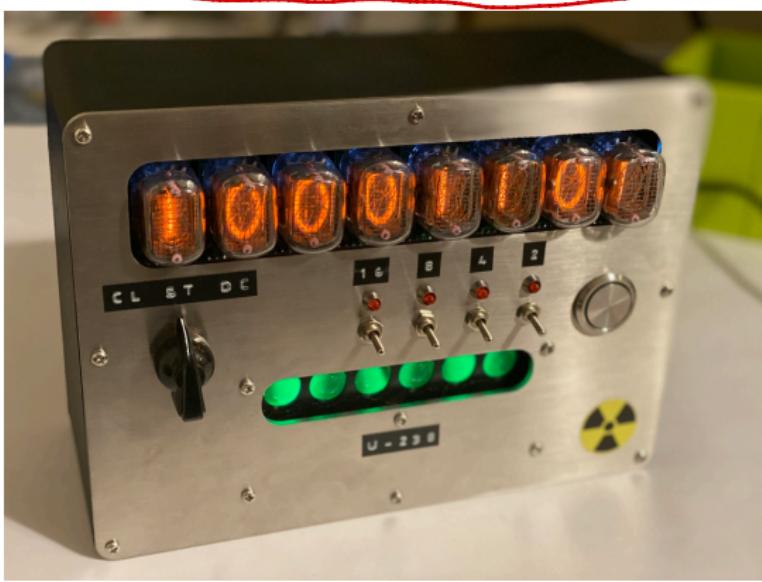


# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

# Randomness Generation

- Earlier days: manual techniques including coin flipping, dice rolling, card shuffling, and roulette spinning
- Early days: physical devices including noise diodes and Geiger counters ([https://github.com/nategri/chernobyl\\_dice](https://github.com/nategri/chernobyl_dice))



# Randomness Generation

- The prevailing belief: only mechanical or electronic devices could produce truly random sequences
- The book: A Million Random Digits With 100,000 Normal Deviates (based on Uranium radiation) RAND
- Current days: computer simulation with deterministic algorithms, also called pseudorandom number generator

Unif(0,1)

# Sampling

- Assuming an algorithm is available for generating  $\text{Unif}(0, 1)$  random numbers
- Two elementary methods for generating random variables (or samples)
  - ▶ Inverse-transform method: operates on the CDF
  - ▶ The acceptance-rejection method: operates on the PDF (or PMF)

# Inverse Transform Method

- Given a  $\text{Unif}(0, 1)$  r.v., we can construct an r.v. with any continuous distribution we want.
- Conversely, given an r.v. with an arbitrary continuous distribution, we can create a  $\text{Unif}(0, 1)$  r.v.
- Other names:
  - ▶ probability integral transform
  - ▶ inverse transform sampling
  - ▶ the quantile transformation
  - ▶ the fundamental theorem of simulation

# Inverse Transform Method: Recall

## Theorem

Let  $F$  be a CDF which is a continuous function and strictly increasing on the support of the distribution. This ensures that the inverse function  $F^{-1}$  exists, as a function from  $(0, 1)$  to  $\mathbb{R}$ . We then have the following results.

- ① Let  $U \sim \text{Unif}(0, 1)$  and  $X = F^{-1}(U)$ . Then  $X$  is an r.v. with CDF  $F$ .
- ② Let  $X$  be an r.v. with CDF  $F$ . Then  $F(X) \sim \text{Unif}(0, 1)$ .

---

## Algorithm Inverse-Transform Method: PDF Case

---

**input:** Cumulative distribution function  $F$ .

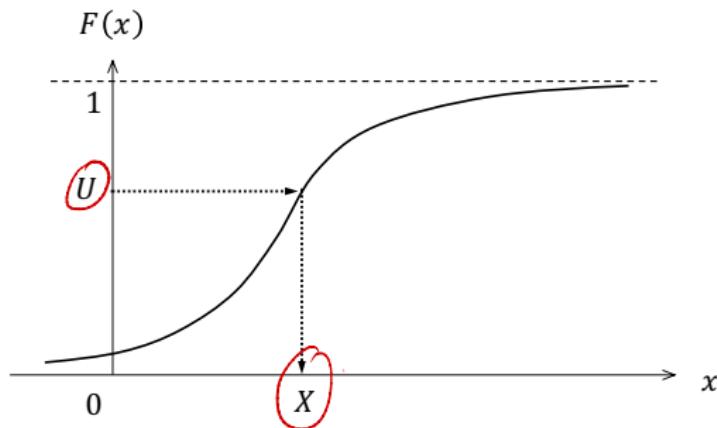
**output:** Random variable  $X$  distributed according to  $F$ .

1: Generate  $U$  from  $\text{Unif}(0, 1)$ .

2:  $X \leftarrow F^{-1}(U)$

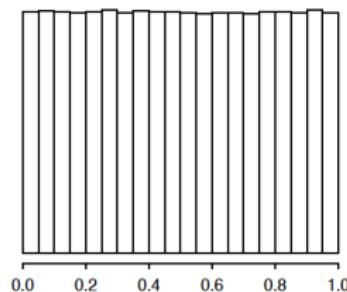
3: **return**  $X$

---

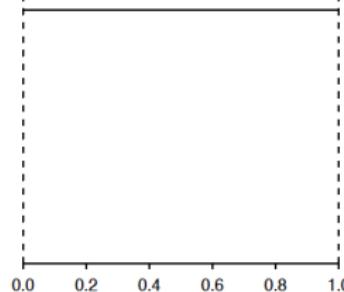


# Histogram & PDF: Example

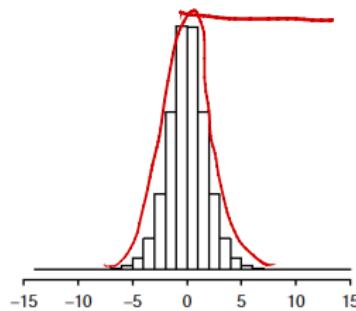
Histogram of  $u$



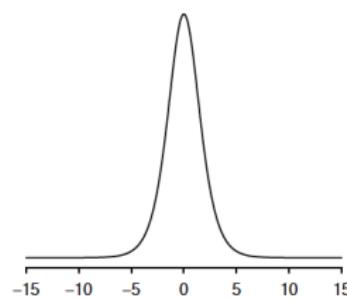
Unif(0,1) PDF



Histogram of  $\log(u/(1-u))$



Logistic PDF



## Box-Muller Method: Recall

$$1^{\circ}. U \sim \text{Unif}(0, 2\pi) = 2\pi \text{ Unif}(0, 1)$$

$$U_2 \sim \text{Unif}(0, 1), \quad U = \underline{2\pi U_2}$$

Let  $U \sim \text{Unif}(0, 2\pi)$ , and let  $T \sim \text{Expo}(1)$  be independent of  $U$ .

Define  $X = \sqrt{2T} \cos U$  and  $Y = \sqrt{2T} \sin U$ . Then  $X$  and  $Y$  are independent, and their marginal distributions are standard normal distribution.

$$2^{\circ}. F_T(t) = 1 - e^{-t}, t > 0 \Rightarrow F_T^{-1}(u) = -\ln(1-u)$$

---

**Algorithm** Normal Random Variable Generation: Box-Muller Approach

$$U_1 \sim \text{Unif}(0, 1), \quad -\ln(1-U_1) \sim \text{Expo}(1)$$

**output:** Independent standard normal random variables  $X$  and  $Y$ .

- 1: Generate two independent random variables,  $U_1$  and  $U_2$ , from  $\text{Unif}(0, 1)$ .
- 2:  $X \leftarrow (-2 \ln U_1)^{1/2} \cos(2\pi U_2)$
- 3:  $Y \leftarrow (-2 \ln U_1)^{1/2} \sin(2\pi U_2)$
- 4: **return**  $X, Y$

$$\sqrt{-2 \ln U_1}$$

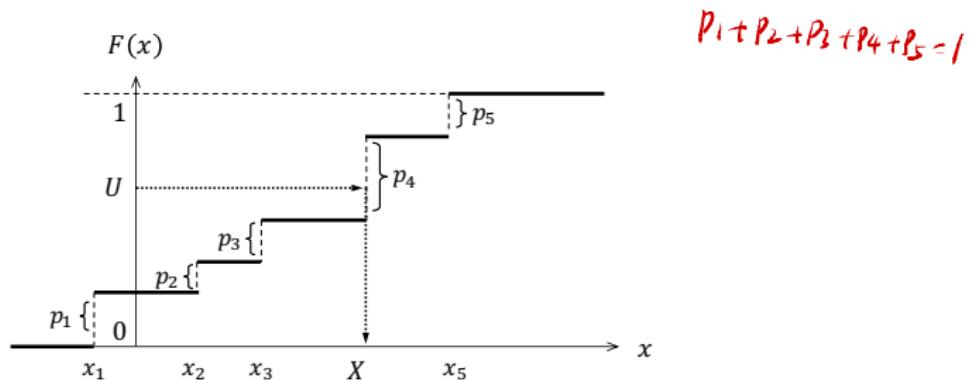
$$U_1 \sim 1 - U_1$$

$$T := -\ln U_1$$

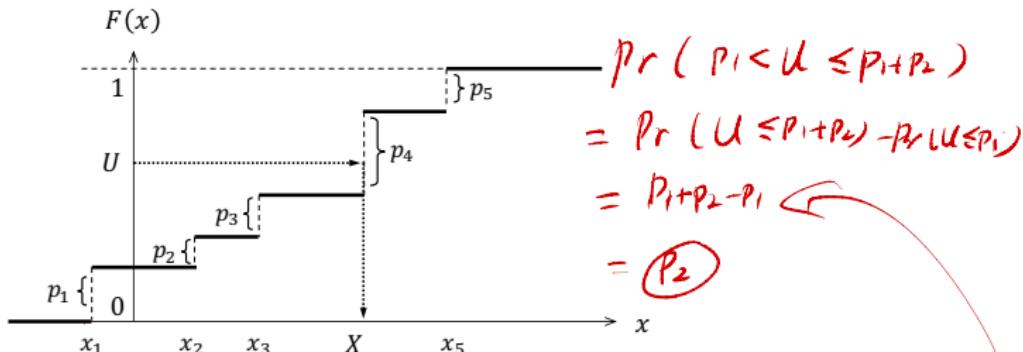
# Inverse Transform Method for Discrete Distribution

- Example: PMF with  $P(X = x_j) = p_j, j = 1, 2, 3, 4, 5$ ,  
 $x_1 < x_2 < x_3 < x_4 < x_5$ .
- CDF with

$$F(x_k) = P(X \leq x_k) = \sum_{j=1}^k P(X = x_j) = \sum_{j=1}^k p_j, k = 1, 2, 3, 4, 5.$$



# Inverse Transform Method for Discrete Distribution



- $U \sim \text{Unif}(0, 1)$ :

$$\Pr(0 < U \leq p_1) = \Pr(U \leq p_1) = (p_1) = \underline{p(X=x_1)}$$

$$X = \begin{cases} x_1 & \text{if } 0 < U \leq p_1 \\ x_2 & \text{if } p_1 < U \leq p_1 + p_2 \\ x_3 & \text{if } p_1 + p_2 < U \leq p_1 + p_2 + p_3 \\ x_4 & \text{if } p_1 + p_2 + p_3 < U \leq p_1 + p_2 + p_3 + p_4 \\ x_5 & \text{if } p_1 + p_2 + p_3 + p_4 < U \leq 1 \end{cases}$$

# Inverse Transform Method for Discrete Distribution

Generalization  
of Inverse Transform

$$F^{-1}(u) = \inf \{x : F(x) \geq u\}$$

## Algorithm Inverse-Transform Method: PMF Case

**input:** Discrete cumulative distribution function  $F$  with monotonic sequence  $\{x_j\}$

**output:** Discrete random variable  $X$  distributed according to  $F$ .

1: Generate  $U \sim \text{Unif}(0, 1)$ .

2: Find the smallest positive integer,  $k$ , such that  $U \leq F(x_k)$ . Let  $X \leftarrow x_k$ .

3: return  $X$

$$\underline{k-1} \quad \underline{U > F(x_{k-1})}$$

$$\begin{aligned} & \Pr(F(x_{k-1}) < U \leq F(x_k)) \\ &= \Pr(U \leq F(x_k)) - \Pr(U \leq F(x_{k-1})) \\ &= F(x_k) - F(x_{k-1}) = \sum_{j=1}^k p_j - \sum_{j=1}^{k-1} p_j = p_k \end{aligned}$$

# Bernoulli Distribution

$$x_1 = 0, x_2 = 1$$

- Bernoulli distribution  $\text{Bern}(p)$  with PMF:  
 $P(X = 1) = p, P(X = 0) = 1 - p, 0 < p < 1.$

$$\begin{array}{l} F(x_1) = 1-p, \\ \text{or} \\ F(x_2) = 1, \end{array}$$

---

## Algorithm Inverse-Transform Method: PMF Case

---

**input:**  $p \in (0, 1)$

**output:** Discrete random variable  $X \sim \text{Bern}(p)$

- 1: Generate  $U \sim \text{Unif}(0, 1)$ .
  - 2: If  $U \leq 1 - p$ , then  $\underline{X \leftarrow 0}$
  - 3: Else  $\underline{X \leftarrow 1}$
  - 4: **return**  $X$
-

# Inverse Transform Method for Discrete Distribution

---

**Algorithm** Inverse-Transform Method: PMF Case

**input:** PMF  $\{p_j\}$  for distribution with non-monotonic sequence  
 $\{x_j\}$

**output:** Discrete random variable  $X$  distributed according to PMF  $\{p_j\}$ .

- 1: Generate  $U \sim \text{Unif}(0, 1)$ .
- 2: Find the positive integer,  $k$ , such that

$$\sum_{j=1}^{k-1} p_j < U \leq \sum_{j=1}^k p_j.$$

$k$

- 3: **return**  $X = x_k$
-

# Bernoulli Distribution

$$X_1 = 1, X_2 = 0$$

- Bernoulli distribution  $\text{Bern}(p)$  with PMF:

$$P(X = 1) = p, P(X = 0) = 1 - p, 0 < p < 1.$$

---

## Algorithm Inverse-Transform Method: PMF Case

---

**input:**  $p \in (0, 1)$

**output:** Discrete random variable  $X \sim \text{Bern}(p)$

1: Generate  $U \sim \text{Unif}(0, 1)$ .

2: If  $U \leq p$ , then  $X \leftarrow 1$

3: Else  $X \leftarrow 0$

4: **return**  $X$

---

# Acceptance-Rejection Method

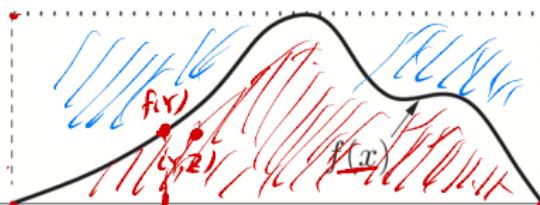
③ Acceptance-Rejection: if  $Z \leq f(Y)$ ; Accept  $(Y, Z)$ .

Red A

$$= \{(y, z) : a \leq y \leq b, 0 \leq z \leq f(y)\}$$

④  $(Y^*, Z^*) \in \text{Red A}$

$$f_{Y^*, Z^*}(y, z) = \frac{1}{\text{Area}(\text{Red})} = 1 \quad Y \Rightarrow f_{Y^*(y)} = \int_0^b f_{Y^*, Z^*}(y, z) dz$$



① Objective PDF  $f: X \in [a, b]$ ,  $\text{C} \geq \sup_x f(x)$

$$\int_a^b f(x) dx = 1 \quad \text{Area}(\text{Red}) =$$

②  $Y \sim \text{Unif}(a, b)$

$Z \sim \text{Unif}(0, c)$  independent

$(Y, Z)$  uniform over the Rectangle.

## Algorithm Acceptance-Rejection Algorithm

Step 1: Generate  $Y \sim \text{Unif}(a, b)$ .

$$= \int_0^{f(y)} 1 \cdot dz = f(y)$$

Step 2: Generate  $Z \sim \text{Unif}(0, c)$ .

$$\Rightarrow Y^* \sim f$$

Step 3: If  $Z \leq f(Y)$ , set  $X = Y$ . Otherwise go back to step 1.

# Acceptance-Rejection Method

②  $Y \sim g$ ;  $Z \sim \text{unif}(0, c \cdot g(Y))$ .

$$\Rightarrow f_{Y,Z}(y,z) = f_Y(y) \cdot f_{Z|Y}(z|y)$$

$$= g(y) \cdot \frac{1}{c \cdot g(y)} = \frac{1}{c}$$

$$= \frac{1}{\text{Area(Triangle)}}.$$

$$\Rightarrow (Y, Z)$$

$$\sim \text{unif}(\text{Triangle})$$

$$\textcircled{3} \quad (Y^*, Z^*) \in \text{Red } A. \quad \sim \text{unif}(A).$$

①  $X \in [a, b]$ , PDF  $f$ : desired

$$\textcircled{2}: \phi(x) = c \cdot g(x) \geq f(x),$$

$$\Rightarrow c \geq \frac{f(x)}{g(x)}, \forall x \in [a, b]$$

$$\Rightarrow c = \sup_x \frac{f(x)}{g(x)}.$$

Area(Triangle)

$$= \int_a^b \phi(x) dx$$

$$\Rightarrow x = \int_a^b c \cdot g(x) dx$$

$$= c \int_a^b g(x) dx$$

$$= c$$

## Algorithm Acceptance-Rejection Algorithm

Step 1: Generate  $Y \sim g$ .

$Z | Y=y \sim \text{unif}(0, c \cdot g(y))$

Step 2: Generate  $Z \sim \text{Unif}(0, c \cdot g(Y))$ .

Step 3: If  $\boxed{Z \leq f(Y)}$ , set  $X = Y$ . Otherwise go back to step 1.

$$\text{Unif}(0, c \cdot g(Y)) \leq f(Y) \Leftrightarrow c \cdot g(Y) \cdot \text{Unif}(0, 1) \leq f(Y)$$

$$\Leftrightarrow \text{Unif}(0, 1) \leq \frac{f(Y)}{c \cdot g(Y)}$$

# Acceptance-Rejection Method

$f = g$   
Supporting set

- Suppose one can generate samples (relatively easily) from PDF  $g$
- How can random samples be simulated from PDF  $f$ ?

---

## **Algorithm** Acceptance-Rejection Algorithm

---

Let  $c$  denote a constant such that  $c \geq \sup_y \frac{f(y)}{g(y)}$ . Then:

Step 1: Generate  $Y \sim g$ .

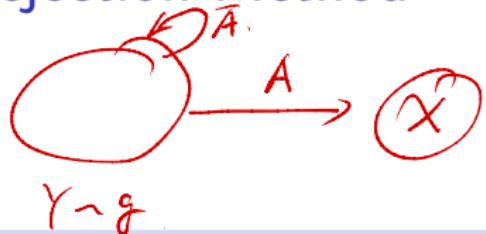
Step 2: Generate  $U \sim \text{Unif}(0, 1)$ .

Step 3: If  $U \leq \frac{f(Y)}{c \cdot g(Y)}$ , set  $X = Y$ . Otherwise go back to step 1.

---

# Acceptance-Rejection Method

(ii)



# of iterations

$$N \sim F_S(p)$$

$$P = P(A) = \frac{1}{c}$$

$$\Rightarrow E(N) = \frac{1}{p} = c$$

## Theorem

- (i) The random variable generated by the Acceptance-Rejection method has the desired PDF f.
- (ii) The number of iterations of the algorithm that are needed is a first-success random variable with mean c.
- (iii)  $c \geq 1$  ✓

Proof (i). event  $A = "U \leq \frac{f(Y)}{C \cdot g(Y)}"$ ,  $f_{Y|A}(y|A) = \underline{\underline{f(y)}}_{\text{desired pdf}}$

$$f_{Y|A}(y|A) = \frac{P(A|Y=y)}{P(A)} f_Y(y).$$

unnormalized

$$1^{\circ} P(A|Y=y) = P(U \leq \frac{f(y)}{C \cdot g(y)} | Y=y) = P(U \leq \frac{f(y)}{C \cdot g(y)})$$

$$= P(U \leq \frac{f(y)}{C \cdot g(y)}) = \frac{f(y)}{C \cdot g(y)}$$

(C ≥ sup\_y  $\frac{f(y)}{g(y)}$ )  
⇒  $\frac{f(y)}{C \cdot g(y)} \leq 1$ .

$$2^{\circ} P(A) \stackrel{\text{LopP}}{\approx} \int P(A|Y=y) \cdot g(y) dy = \int \frac{f(y)}{C \cdot g(y)} g(y) dy$$

$$= \frac{1}{C} \int f(y) dy = \frac{1}{C} \cdot \leq 1 \Rightarrow C \geq 1.$$

$$\Rightarrow f_{Y|A}(y|A) = \frac{P(A|Y=y)}{P(A)} \cdot \underline{\underline{f_Y(y)}} = \frac{\frac{f(y)}{C \cdot g(y)}}{\frac{1}{C}} \cdot g(y) = f(y).$$

$\Rightarrow X \sim f$ .

# Proof

## Example: Beta Distribution

$$X \sim \text{Beta}(2, 4), \quad f(x) = 20x(1-x)^3, \quad 0 < x < 1$$

- An r.v.  $X$  is said to have the *Beta distribution* with parameters  $a$  and  $b$ ,  $a > 0$  and  $b > 0$ , if its PDF is

$$f(x) = \frac{1}{\beta(a, b)} x^{a-1} (1-x)^{b-1}, \quad 0 < x < 1,$$

---

where the constant  $\beta(a, b)$  is chosen to make the PDF integrate to 1. We write this as  $X \sim \text{Beta}(a, b)$ .  $a=b=1, f(x) \propto \text{constant}$ .

- Beta distribution is a generalization of uniform distribution.
- Use the Acceptance-Rejection Method to generate a random variable with distribution Beta(2, 4)

# Solution

$$f(x) = 20x(1-x)^3, 0 < x < 1$$

①  $g : \text{Unif}(0,1)$ ,  $g(x) = 1, 0 < x < 1$ .

$$C \geq \sup_{y \in [0,1]} \frac{f(y)}{g(y)} = \sup_{y \in [0,1]} \frac{20y(1-y)^3}{1} = \sup_{y \in [0,1]} 20y(1-y)^3. \Rightarrow y^* = \frac{1}{4}$$
$$\Rightarrow C \geq \frac{135}{64}. \text{ choose } C = \frac{135}{64}.$$

②  $\Rightarrow 0 < y < 1$ ,  $\frac{f(y)}{C \cdot g(y)} = \frac{20y(1-y)^3}{\frac{135}{64} \cdot 1} = \frac{256}{27} y(1-y)^3$ .

---

Step 1 : Generate  $Y \sim \text{Unif}(0,1)$ .

Step 2 : Generate  $U \sim \text{unif}(0,1)$

Step 3 : If  $U \leq \frac{f(Y)}{C \cdot g(Y)} = \frac{256}{27} Y(1-Y)^3$ , set  $X=Y$ .

otherwise reject  $Y$ , go back to step 1.

# Solution

# Example: Normal Distribution

①  $Z \sim N(0,1)$ .  $(-\infty, +\infty)$

$$\underline{X=|Z|} \quad (0, +\infty)$$

$$P(X \leq x) = P(|Z| \leq x) = 2P(0 \leq Z \leq x) = 2 \int_0^x \frac{1}{\sqrt{\pi}} e^{-z^2} dz$$

$$= \int_0^x \sqrt{\frac{2}{\pi}} e^{-z^2} dz \Rightarrow f_X(x) = \sqrt{\frac{2}{\pi}} e^{-\frac{1}{2}x^2}, 0 < x < \infty$$

② choose  $g \sim \text{Exp}(1)$  .  $g(x) = e^{-x}, 0 < x < \infty$ .

- Use the Acceptance-Rejection Method to generate a random variable with distribution  $N(0, 1)$

$$C \geq \sup_y \frac{f(y)}{g(y)} = \sup_y \sqrt{\frac{2}{\pi}} e^{-\frac{1}{2}y^2} / e^{-y} = \sup_y \sqrt{\frac{2}{\pi}} e^{-\frac{1}{2}(y-1)^2 + \frac{1}{2}} = \sqrt{\frac{2e}{\pi}}$$

$$\Rightarrow y^* = 1, \text{ choose } C = \sqrt{\frac{2e}{\pi}}$$

$$\Rightarrow \frac{f(y)}{C \cdot g(y)} = e^{1(y - \frac{1}{2}y^2 - \frac{1}{2})} = e^{-\frac{1}{2}(y-1)^2}$$

### Solution ③

Step 1 : Generate  $Y \sim \text{Exp}(1)$

2 :  $\dots \sim U \sim \text{unif}(0,1)$

3 : If  $U \leq e^{-\frac{1}{2}(Y-1)^2}$ , set  $X = Y$ .

otherwise return to Step 1.

$X \sim \mathcal{N}(0,1)$

Step 4 : generate  $u' \sim \text{unif}(0,1)$

Box-Muller

Acceptance-Rejection

$$Z = \begin{cases} X & \text{if } u' \leq \frac{1}{2} \\ -X & \text{otherwise.} \end{cases}$$

$$\underline{Z \sim \mathcal{N}(0,1)}$$

# Solution

# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation  $f(x,y,z) = \underline{f(x)} \cdot \underline{f(y|x)} \cdot \underline{f(z|x,y)}$
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

# Change of Variables

## Theorem

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a continuous random vector with joint PDF  $f_{\mathbf{X}}(x)$ , and let  $\mathbf{Y} = g(\mathbf{X})$  where  $g$  is an invertible function from  $\mathbb{R}^n$  to  $\mathbb{R}^n$ . Let  $y = g(\mathbf{x})$  and suppose that all the partial derivatives  $\frac{\partial x_i}{\partial y_j}$  exists and are continuous, so we can form the **Jacobian matrix**

$$\frac{\partial \mathbf{x}}{\partial \mathbf{y}} = \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \cdots & \frac{\partial x_1}{\partial y_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial x_n}{\partial y_1} & \frac{\partial x_n}{\partial y_2} & \cdots & \frac{\partial x_n}{\partial y_n} \end{pmatrix}$$

Also assume that the determinant of the Jacobian matrix is never 0. Then the joint PDF of  $\mathbf{Y}$  is

$$f_{\mathbf{Y}}(y) = f_{\mathbf{X}}(x) \left| \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right|$$

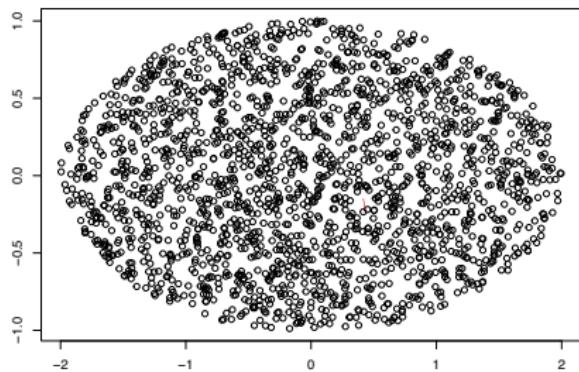
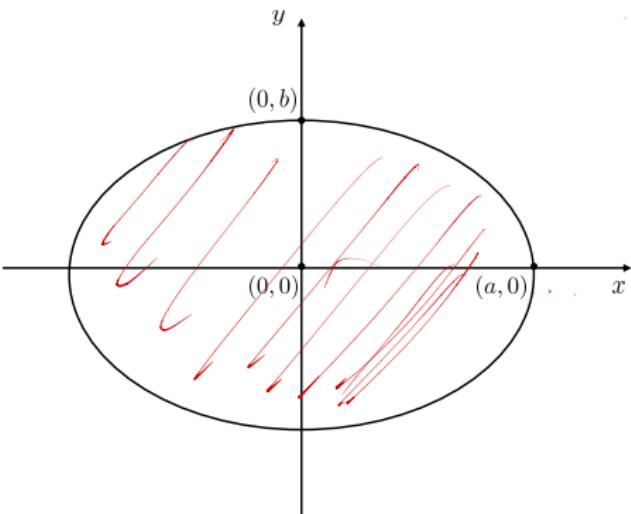
# Example: Generate Uniform Distribution over An Ellipse

objective PDF  $f_{X,Y}(x,y) = \frac{1}{\pi \cdot a \cdot b} \quad \theta(x,y) \in E_2(a,b)$

- Ellipse:

$$x = \rho \cdot a \cdot \cos \theta \quad \rho \in [0,1] ; \theta \in [0,2\pi]$$

$$E_2(a,b) = \left\{ (x,y) \in \mathbb{R}^2 : \left(\frac{x}{a}\right)^2 + \left(\frac{y}{b}\right)^2 \leq 1 \right\}$$



Solution ① Jacobi Matrix

$$\begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{bmatrix} = \begin{bmatrix} a \cos \theta & -r a \sin \theta \\ b \sin \theta & r b \cos \theta \end{bmatrix}$$

$$J = \det(J) \checkmark, J = r a b.$$

$$\Rightarrow f_{R,\Theta}(r, \theta) = f_{x,y}(x, y) \cdot |J| = \frac{1}{\pi a b} r a b = \frac{1}{\pi}, r \in [0, 1], \theta \in [0, \pi].$$

$$\Rightarrow f_R(r) = \int_0^{2\pi} \frac{1}{\pi} d\theta = 2\pi, 0 \leq r \leq 1 \Rightarrow F_R(r) = \underline{r^2}, 0 \leq r \leq 1.$$

$$f_\Theta(\theta) = \int_0^{2\pi} \frac{1}{\pi} d\theta = \frac{1}{2\pi}, 0 \leq \theta \leq 2\pi. \quad \theta \sim \text{Unif}(0, 2\pi)$$

$$\Rightarrow f_{R,\Theta}(r, \theta) = f_R(r) \cdot f_\Theta(\theta), R, \Theta \text{ independent.}$$

(R)  $\cdot F_R^{-1}(z) = \sqrt{z}, 0 \leq z \leq 1;$

Solution ②

$U_1, U_2$  independent Uniform

$$X = a \sqrt{U_1} \cos(2\pi U_2)$$

$$Y = b \sqrt{U_1} \sin(2\pi U_2)$$

$$\underline{(X, Y)}$$

# Solution

# Change of Variables

$A : n \times n$

$$\det(ATA) = \det^2(A)$$

## Theorem

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a continuous random vector with joint PDF  $f_{\mathbf{X}}(x)$ , and let  $\mathbf{Y} = g(\mathbf{X})$  where  $g$  is a function from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ . Let  $y = g(x)$  and we have the Jacobian matrix  $\frac{\partial \mathbf{x}}{\partial \mathbf{y}}$ . The corresponding Gram matrix is

$$\mathbf{G} = \left( \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right)^T \frac{\partial \mathbf{x}}{\partial \mathbf{y}}.$$

Then the joint PDF of  $\mathbf{Y}$  is

$$f_{\mathbf{Y}}(y) = f_{\mathbf{X}}(x) \sqrt{\det(\mathbf{G})}$$

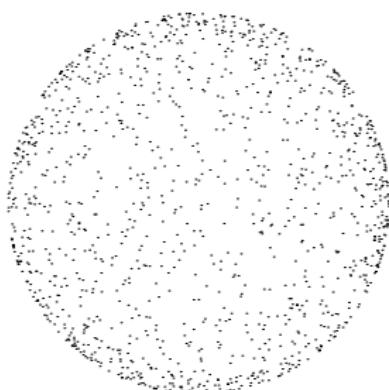
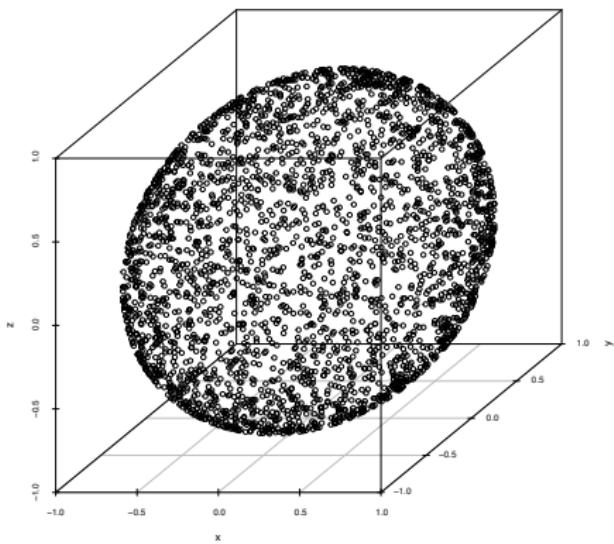
# Example: Generate Uniform Distribution over A

Sphere

$$\text{Ball } B_3(r) = \{(x,y,z) \in \mathbb{R}^3 : x^2 + y^2 + z^2 \leq r^2\}$$

- Sphere:

$$S_2(r) = \{(x, y, z) \in \mathbb{R}^3 : x^2 + y^2 + z^2 = r^2\}.$$



Solution ①  $f_{x,y,z}(x,y,z) = \frac{1}{4\pi r^2}$  ( $(x,y,z) \in S_{out}$ ).

$$\begin{cases} x = r \sin \theta \cos \phi \\ y = r \sin \theta \sin \phi \\ z = r \cos \theta \end{cases} \quad \begin{array}{l} \theta \in [0, \pi], \\ \phi \in [0, 2\pi], \end{array} \quad \underline{(x,y,z)} \rightarrow \underline{(\theta, \phi)}$$

Jacobi Matrix  $M = \begin{bmatrix} \frac{\partial x}{\partial \theta} & \frac{\partial x}{\partial \phi} \\ \frac{\partial y}{\partial \theta} & \frac{\partial y}{\partial \phi} \\ \frac{\partial z}{\partial \theta} & \frac{\partial z}{\partial \phi} \end{bmatrix} = \begin{bmatrix} r \cos \theta \cos \phi & -r \sin \theta \sin \phi \\ r \cos \theta \sin \phi & r \sin \theta \cos \phi \\ -r \sin \theta & 0 \end{bmatrix}$

Gram matrix  $G = M^T M = \begin{bmatrix} r^2 & 0 \\ 0 & r^2 \sin^2 \theta \end{bmatrix} \Rightarrow \det(G) = r^4 \sin^2 \theta.$

$$\Rightarrow f_{\theta, \phi}(\theta, \phi) = f_{x,y,z}(x,y,z) \cdot \sqrt{\det(G)} = \frac{1}{4\pi r^2} \cdot r^2 \sin \theta = \frac{1}{4\pi} \sin \theta$$

$\theta \in [0, \pi]$   
 $\phi \in [0, 2\pi]$

Solution ②  $f_{\theta}(\theta) = \int_0^{2\pi} f_{\theta,\phi}(\theta, \phi) d\phi = \frac{1}{2} \sin \theta ; 0 \leq \theta \leq \pi.$

$$F_{\theta}(\theta) = \int_0^{\theta} f_{\theta}(s) ds = \left( \frac{1-\cos \theta}{2} \right), 0 \leq \theta \leq \pi.$$

$$f_{\phi}(\phi) = \frac{1}{2\pi} ; 0 \leq \phi \leq 2\pi ; \Rightarrow f_{\theta,\phi}(\theta, \phi) = f_{\theta}(\theta) f_{\phi}(\phi)$$

$\theta$  and  $\phi$  independent. ;  $\phi \sim \text{Unif}(0, 2\pi) = 2\pi \text{ Unif}(0, 1);$

$$\underline{F_{\theta}^{-1}(s) = \arccos(1-2s)}, 0 \leq s \leq 1$$

Note: if  $\theta = \arccos(1-2s) \Rightarrow \cos \theta = 1-2s$

$$\sin^2 \theta = 1 - \cos^2 \theta = 1 - (1-2s)^2 = 4s(1-s)$$

$$\Rightarrow \sin \theta > 0 ; \Rightarrow \underline{\sin \theta = 2\sqrt{s(1-s)}}$$

### Solution ③

independently generate  $U_1, U_2 \sim \text{Uniform}$

$$\left\{ \begin{array}{l} X \leftarrow r \cdot 2\sqrt{U_1(1-U_1)} \cdot \cos(2\pi U_2) \\ Y \leftarrow r \cdot 2\sqrt{U_1(1-U_1)} \cdot \sin(2\pi U_2) \\ Z \leftarrow r \cdot (1-2U_1) \end{array} \right.$$

# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

# Monte Carlo Integration

~~E(X)~~

$$\underline{x}_1, \dots, \underline{x}_n : \frac{1}{n} (x_1 + \dots + x_n)$$
$$E[g(x)] \quad g(x_1), \dots, g(x_n) : \frac{1}{n} (g(x_1) + \dots + g(x_n))$$

- We can use the sample mean to approximate the expectation:

$$\underline{E[g(X)]} \approx \frac{1}{n} \sum_{i=1}^n g(X_i). \quad X \sim \text{Unif}(a, b).$$

- Now we have integration

$$\int_a^b g(x) dx = (b - a) \int_a^b g(x) \cdot \frac{1}{b - a} dx. \quad f(x) \cdot \text{PDF}$$
$$= (b - a) \int_a^b g(x) dx = (b - a) E[g(X)]$$

- Drawing n samples (empirical samples) from Unif(a, b):

$$\underline{X_1, X_2, \dots, X_n \sim \text{Unif}(a, b)}.$$

$$\pi(b-a) \cdot \frac{1}{n} (g(x_1) + \dots + g(x_n))$$

- Monte Carlo Integration:

$$\int_a^b g(x) dx \approx \frac{1}{n} \sum_{i=1}^n g(X_i) (b - a).$$

# Monte Carlo Integration

# Example: $\pi$ as An Integration

Evaluate the integration

$$\int_0^1 \frac{4}{1+x^2} dx.$$

- $g(x) = 4/(1+x^2)$ ,  $0 < x < 1$ .  $a=0, b=1$
- $X_1, \dots, X_n$ : samples from  $\text{Unif}(0, 1)$ .
- Monte Carlo Integration:

$$\int_0^1 \frac{4}{1+x^2} dx \approx \frac{1}{n} \sum_{i=1}^n \left( \frac{4}{1+X_i^2} \right)$$

# Example

Evaluate the integration

$$\int_0^4 \sqrt{x + \sqrt{x + \sqrt{x + \sqrt{x}}}} dx.$$

- Corresponding

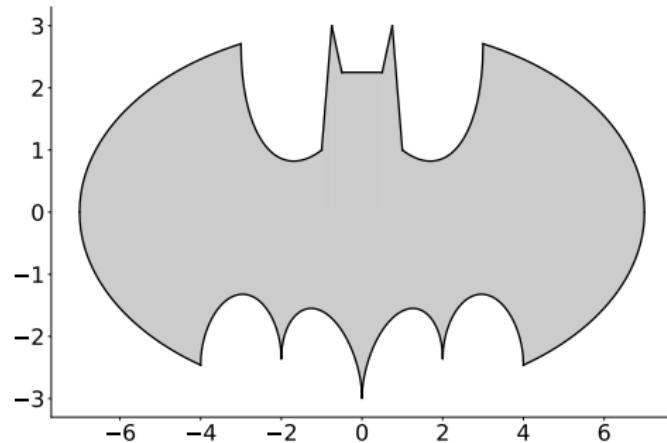
$$g(x) = \sqrt{x + \sqrt{x + \sqrt{x + \sqrt{x}}}}$$

- $X_1, \dots, X_n$ : samples from  $\text{Unif}(0, 4)$ .  $a=0, b=4$
- Monte Carlo Integration:

$$\int_0^4 \sqrt{x + \sqrt{x + \sqrt{x + \sqrt{x}}}} dx \approx \frac{4}{n} \sum_{i=1}^n \sqrt{X_i + \sqrt{X_i + \sqrt{X_i + \sqrt{X_i}}}}$$

# Example: Area of Batman Curve

- Challenging and Fun
- <https://mathworld.wolfram.com/BatmanCurve.html>



# Example: Estimation of Probability

- Indicator: bridge between expectation and probability
- Given event  $A$ :

$$\underline{I_A(x)} = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{Otherwise} \end{cases}.$$

- For random variable  $X$ :

$$\underline{P(X \in A)} = 1 \cdot P(X \in A) + 0 \cdot P(X \notin A)$$
$$= \underline{E(I_A(X))}$$

$$\approx \frac{1}{n} \sum_{i=1}^n \underline{I_A(X_i)}.$$

$x_1, \dots, x_n \sim X$

## Example: Estimation of $\pi$

generate  $n$  points  $(X_1, Y_1), \dots, (X_n, Y_n)$

① event  $A_i$ : "the  $i^{th}$  point lands within the circle",  $\Leftrightarrow \{X_i^2 + Y_i^2 \leq 1\}$

$$-1 \leq X_i \leq 1$$

$$-1 \leq Y_i \leq 1$$

②  $I_{A_i} = Z_i$

$$P(Z_i = 1) = P(A_i) = \frac{\pi \cdot 1}{4} = \frac{\pi}{4}$$

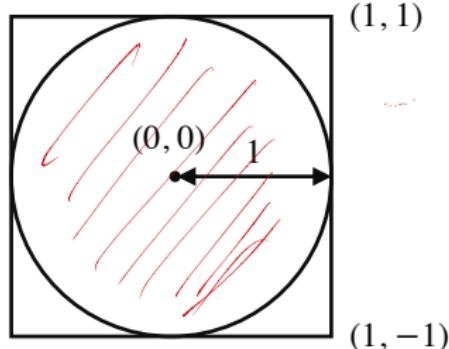
$$P(Z_i = 0) = 1 - \frac{\pi}{4}$$

$$\Rightarrow E(Z_i) = P(Z_i = 1) = \frac{\pi}{4}$$

$$Z_1, \dots, Z_n \sim \underbrace{\text{Bern}(\frac{\pi}{4})}_{(-1, -1)}$$

$$\Rightarrow E(Z) = \frac{\pi}{4}$$

$$\Rightarrow \pi = 4E(Z) \approx 4 \cdot \underbrace{\frac{1}{n}(Z_1 + \dots + Z_n)}_{\text{--}}$$



# Example: Estimation of $\pi$

# Example: Estimation of $\pi$

# Useful Tools: Importance Sampling

- Standard Monte Carlo integration is great if you can sample from the target distribution (i.e. the desired distribution)
- But what if you can't sample from the target?
- **Importance Sampling:** draw the sample from a proposal distribution and re-weight the integral using importance weights so that the correct distribution is targeted

# Importance Sampling

$$H = \underline{E_f[h(Y)]} = \int \underline{h(y)f(y)dy}$$

- $h$  is some function and  $f$  is the PDF of random variable  $Y$
- When the PDF  $f$  is difficult to sample from, importance sampling can be used
- Rather than sampling from  $f$ , you specify a different PDF  $g$ , as the proposal distribution.

$$\underline{H = \int h(y)f(y)dy} = \int h(y) \frac{f(y)}{g(y)} g(y)dy = \int \boxed{\frac{h(y)f(y)}{g(y)}} g(y)dy$$

# Importance Sampling

$g$  is a PDF.

$$H = E_f[h(Y)] = \int \frac{h(y)f(y)}{g(y)} dy = E_g\left[\frac{h(Y)f(Y)}{g(Y)}\right]$$

- Hence, given an iid sample  $Y_1, \dots, Y_n$  from PDF  $\underline{g}$ , our estimator of  $H$  becomes

$$\hat{H} = \frac{1}{n} \sum_{j=1}^n \frac{h(Y_j)f(Y_j)}{g(Y_j)}$$

# Example: Gaussian Tail Probability

$$P(-3 < Y < 3) = 0.997$$

Method 1:  $C = P(Y > 8) = \frac{E[I(Y > 8)]}{f(Y_1, \dots, Y_n)}$

$C \approx 0$

$$\approx \frac{1}{n} \sum_{j=1}^n I(Y_j > 8)$$

$$f(Y_1, \dots, Y_n) \sim N(0, 1)$$

$$h(y) = I(Y > 8) = \begin{cases} 1 & \text{if } y > 8 \\ 0 & \text{otherwise} \end{cases}$$

Evaluate the probability of rare event  $c = \underline{\mathbb{P}(Y > 8)}$ , where  $\underline{Y \sim N(0, 1)}$ .

choose  $g \sim N(8, 1)$ ,  $Y_1, \dots, Y_n \sim g$ .

Method 2:

importance sampling

$$\begin{aligned} C &\approx \frac{1}{n} \sum_{j=1}^n \frac{h(Y_j) f(Y_j)}{g(Y_j)} = \frac{1}{n} \sum_{j=1}^n I(Y_j > 8) \cdot \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} Y_j^2}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (Y_j - 8)^2}} \\ &= \frac{1}{n} \sum_{j=1}^n I(Y_j > 8) \cdot e^{-8Y_j + 32} \end{aligned}$$

$$n = 50000 \rightarrow C \approx 6.025 \times 10^{-16}$$

# Solution

# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

$$n \rightarrow \infty$$

# Sample Mean: Recall

## Definition

Let  $X_1, \dots, X_n$  be i.i.d. random variables with finite mean  $\mu$  and finite variance  $\sigma^2$ . The *sample mean*  $\bar{X}_n$  is defined as follows:

$$\bar{X}_n = \frac{1}{n} \sum_{j=1}^n X_j.$$

The sample mean  $\bar{X}_n$  is itself an r.v. with mean  $\mu$  and variance  $\sigma^2/n$ .

$$n \rightarrow \infty$$

$$\sigma^2/n \rightarrow 0$$

# Strong Law of Large Numbers (SLLN)

$$\int_a^b g(x)dx \approx \frac{b-a}{n} \sum_{i=1}^n g(x_i) \Rightarrow \frac{b-a}{n} \sum_{i=1}^n g(x_i) \xrightarrow{\text{w.p.1}} \int_a^b g(x)dx.$$

$x_1, \dots, x_n$  i.i.d. unif(a,b).

(1)  $x_1, \dots, x_n$  i.i.d.  $g$  = continuous function.

$g(x_1), \dots, g(x_n)$  i.i.d.

$$E[g(x_i)] = \int_a^b g(x) \cdot \frac{1}{b-a} dx.$$

## Theorem

The sample mean  $\bar{X}_n$  converges to the true mean  $\mu$  pointwise as  $n \rightarrow \infty$ , with probability 1. In other words, the event  $\bar{X}_n \rightarrow \mu$  has probability 1.

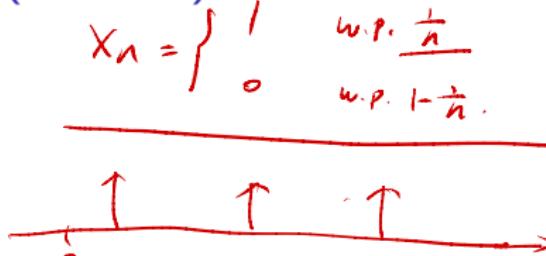
(2) By SLLN  $\frac{g(x_1) + \dots + g(x_n)}{n} \xrightarrow[n \rightarrow \infty]{\text{w.p.1}} E[g(x)] = \int_a^b g(x) \frac{1}{b-a} dx$

$$\Rightarrow \frac{(b-a)}{n} \sum_{i=1}^n g(x_i) \xrightarrow[n \rightarrow \infty]{\text{w.p.1}} \int_a^b g(x)dx.$$

# Weak Law of Large Numbers (WLLN)

$$\begin{aligned} X_n &\xrightarrow{\text{a.s.}} 0 & X_n &\xrightarrow{\text{P}} 0 \\ X_n &\xrightarrow{\text{w.p.1}} 0 & \lim_{n \rightarrow \infty} P(|X_n - 0| > \varepsilon) &= 0 \\ P(\lim_{n \rightarrow \infty} X_n = 0) &= 1 & \forall \varepsilon > 0. \end{aligned}$$

$X_n = \begin{cases} 1 & \text{w.p. } \frac{1}{n} \\ 0 & \text{w.p. } 1 - \frac{1}{n}. \end{cases}$



## Theorem

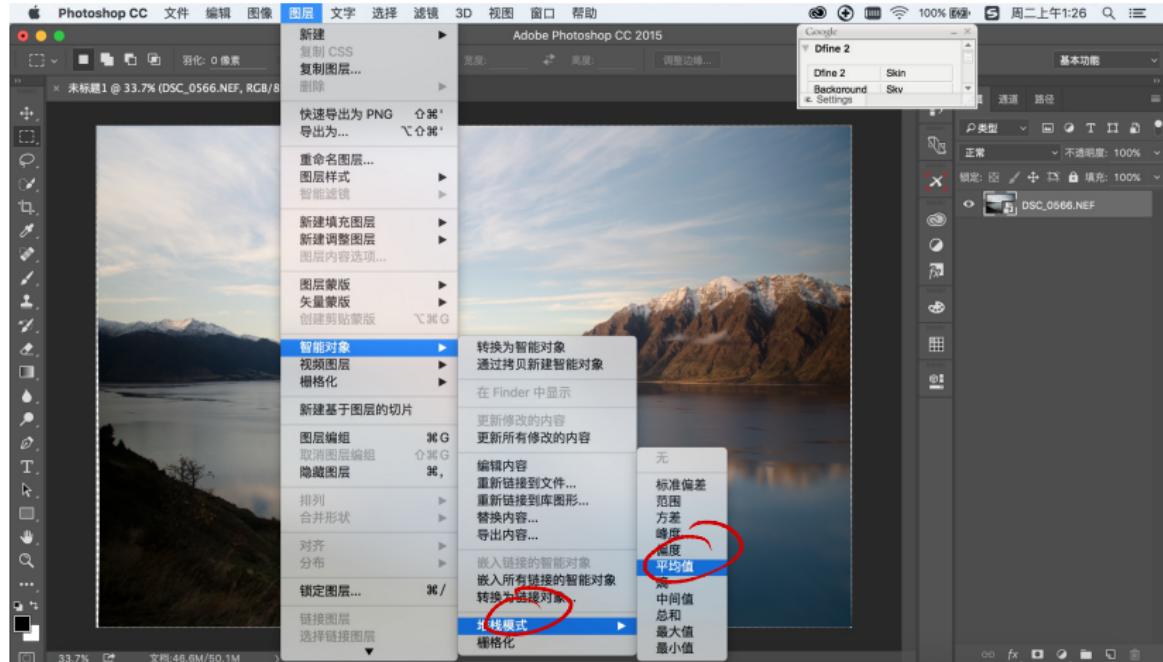
For all  $\epsilon > 0$ ,  $P(|\bar{X}_n - \mu| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$ . (This form of convergence is called convergence in probability).

$$\begin{aligned} \lim_{n \rightarrow \infty} P(|X_n - 0| > \varepsilon) &= P(X_n > \varepsilon) & (\varepsilon > 0) \\ &= \overbrace{P(X_n = 1)}^{0 < \varepsilon < 1} & = 0 \\ &= \frac{1}{n} & \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

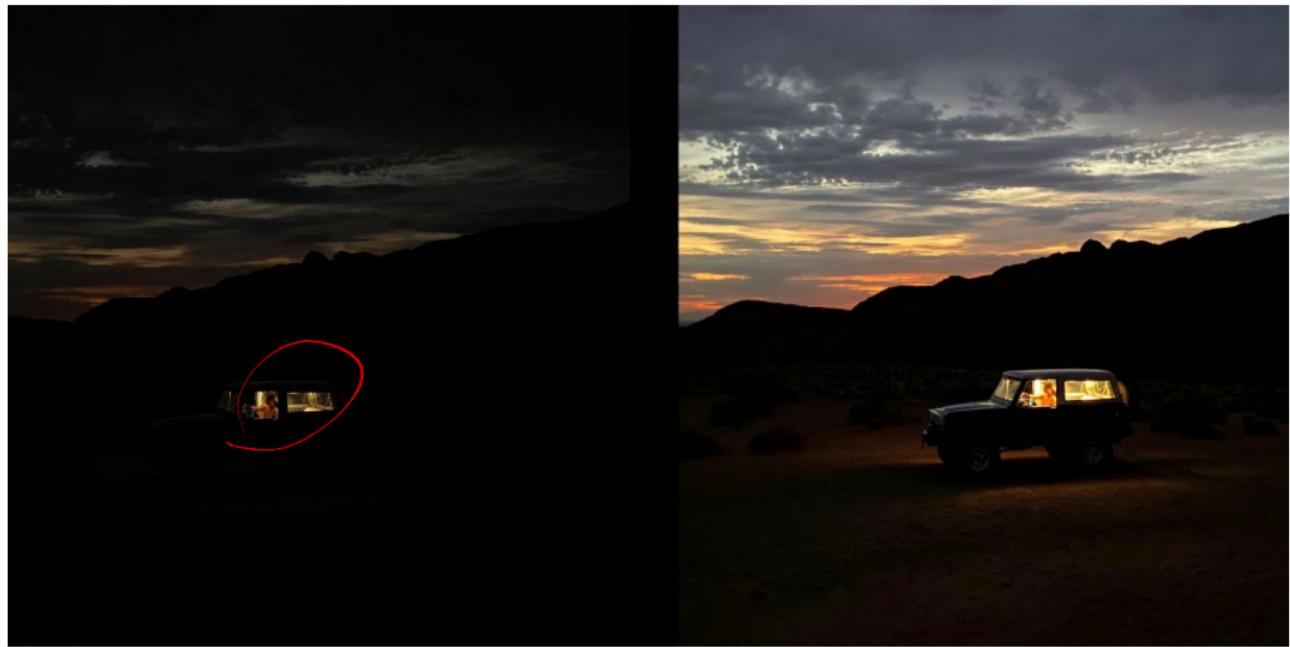
# Widely Applications: Photo Stacking with PC



# Widely Applications: Photo Stacking with PC



# Widely Applications: Night Model with Smart Phone



# Widely Applications: Photo Stacking with Smart Phone



# Widely Applications: Photo Stacking with Smart Phone



# Widely Applications: Photo Stacking with Smart Phone



# Outline

- 1 History of Monte Carlo
- 2 Sampling: Random Variable Generation
- 3 Sampling: Random Vector Generation
- 4 Monte Carlo Integration
- 5 Asymptotic Analysis: Law of Large Numbers
- 6 Non-asymptotic Analysis: Inequalities

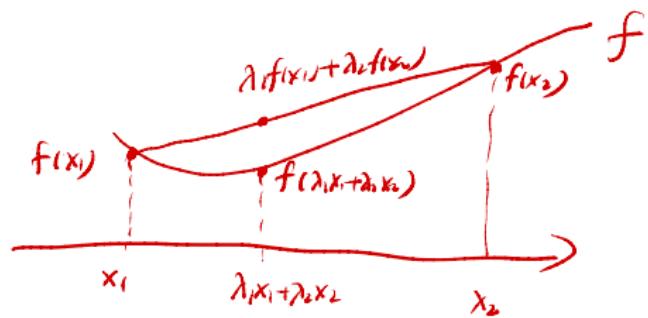
# Cauchy-Schwarz Inequality: Recall

## Theorem

For any r.v.s  $X$  and  $Y$  with finite variances,

$$|E(XY)| \leq \sqrt{E(X^2) E(Y^2)}.$$

## Jensen's Inequality



If  $f$  is a convex function,  $0 \leq \lambda_1, \lambda_2 \leq 1, \lambda_1 + \lambda_2 = 1$ , then for any  $x_1, x_2$ ,

$$f(\lambda_1 x_1 + \lambda_2 x_2) \leq \lambda_1 f(x_1) + \lambda_2 f(x_2).$$

# Jensen's Inequality

## Theorem

Let  $X$  be a random variable. If  $g$  is a convex function, then  $E(g(X)) \geq g(E(X))$ . If  $g$  is a concave function, then  $E(g(X)) \leq g(E(X))$ . In both cases, the only way that equality can hold is if there are constants  $a$  and  $b$  such that  $g(X) = a + bX$  with probability 1.

## Quick Examples

$$\left. \begin{array}{l} g \text{ is convex ; } E[g(X)] \geq g(E(X)) \\ g \text{ is concave ; } E[g(X)] \leq g(E(X)) \end{array} \right\} \begin{array}{l} g'' \geq 0 \\ \text{Convex} \end{array}$$

1<sup>o</sup>.  $g(x) = x^2$ ,  $x \in \mathbb{R}$ ; convex;  $\Rightarrow E[x^2] \geq (E(x))^2$  ✓

$$\text{Var}(x) = E(x^2) - E(x)^2 \geq 0.$$

2<sup>o</sup>.  $g(x) = \frac{1}{x}$ ,  $x > 0$ ; convex;  $\Rightarrow E[\frac{1}{x}] \geq \frac{1}{E(x)}$  ✓

3<sup>o</sup>.  $g(x) = \log x$ ,  $x > 0$ ; concave;  $\Rightarrow E[\log x] \leq \log(E(x))$ .

# Entropy

- Let  $X$  be a discrete r.v. whose distinct possible values are  $a_1, a_2, \dots, a_n$ , with probabilities  $p_1, p_2, \dots, p_n$  respectively (so  $p_1 + p_2 + \dots + p_n = 1$ ).
- The *entropy* of  $X$  is defined as follows:  
$$H(X) = \sum_{j=1}^n p_j \log_2 (1/p_j).$$
- Using Jensen's inequality, show that the maximum possible entropy for  $X$  is when its distribution is uniform over  $a_1, a_2, \dots, a_n$ , i.e.,  $p_j = 1/n$  for all  $j$ .
- This makes sense intuitively, since learning the value of  $X$  conveys the most information on average when  $X$  is equally likely to take any of its values, and the least possible information if  $X$  is a constant.

Proof ① Construct a random variable  $Y$ . s.t

$$Y = \begin{cases} \frac{1}{p_1} & \text{w.p. } p_1 \\ \frac{1}{p_2} & \text{w.p. } p_2 \\ \vdots & \\ \frac{1}{p_n} & \text{w.p. } p_n \end{cases} \Rightarrow E(Y) = \frac{1}{p_1} \cdot p_1 + \frac{1}{p_2} \cdot p_2 + \dots + \frac{1}{p_n} \cdot p_n = n$$

$$\textcircled{2} H(X) \triangleq \sum_{j=1}^n p_j \log_2 \frac{1}{p_j} = \underbrace{E[\log_2 Y]}_{\text{w.p. } p_1, \dots, p_n} \leq \log_2 E[Y] = \log_2 n.$$

$$\text{w.p. } p_1, \dots, p_n \Rightarrow \max_{p_1, \dots, p_n} H(X) \leq \log_2 n$$

$$\textcircled{3} \text{ when } X \sim \text{Dunif}(\frac{1}{n}), p_1 = p_2 = \dots = p_n = \frac{1}{n}; H(X) = \sum_{j=1}^n \frac{1}{n} \cdot \log_2 n = \log_2 n.$$

$$\Rightarrow \max_{p_1, \dots, p_n} H(X) \geq \log_2 n \Rightarrow \max_{p_1, \dots, p_n} H(X) = \log_2 n$$

# Kullback-Leibler Divergence

Let  $\mathbf{p} = (p_1, \dots, p_n)$  and  $\mathbf{r} = (r_1, \dots, r_n)$  be two probability vectors (so each is nonnegative and sums to 1). Think of each as a possible PMF for a random variable whose support consists of  $n$  distinct values. The *Kullback-Leibler* divergence between  $\mathbf{p}$  and  $\mathbf{r}$  is defined as

$$D(\mathbf{p}, \mathbf{r}) = \sum_{j=1}^n p_j \log_2 (1/r_j) - \sum_{j=1}^n p_j \log_2 (1/p_j).$$

Show that the Kullback-Leibler divergence is nonnegative.

Proof ①  $D(P, r) = \sum_{j=1}^n p_j \log_2 \frac{r_j}{p_j} - \sum_{j=1}^n p_j \log_2 \frac{p_j}{p_j} = \sum_{j=1}^n p_j \log_2 \frac{p_j}{r_j}$

 $= -\frac{\sum_{j=1}^n p_j \log_2 \frac{r_j}{p_j}}{\sum_{j=1}^n p_j}$

② Construct a random variable  $Y$ . s.t.

$P(Y = \frac{r_j}{p_j}) = p_j, j = 1, 2, \dots, n.$

$\Rightarrow E(Y) = \sum_{j=1}^n \frac{r_j}{p_j} \cdot p_j = \sum_{j=1}^n r_j = 1$

③  $D(P, r) = -E[\log_2 Y] \geq -\log_2 E(Y) = -\log_2 1 = 0$

# Markov's Inequality

Concentration Inequality

Chebyshev  
Markov      Lapunov

$$P(|X - E(X)| \geq a) \leq \frac{1}{\frac{1}{a^2} e^{-a}}$$

## Theorem

For any r.v.  $X$  and constant  $a > 0$ ,

$$P(|X| \geq a) \leq \frac{E|X|}{a}.$$

# Proof

$$\boxed{P(|X| \geq a) \leq \frac{1}{a} E[|X|], \quad a > 0}$$

①  $Y = \frac{1}{a} |X| \geq 0 : \quad \underline{I(Y \geq 1)} \leq Y \quad \begin{cases} Y \geq 1 & LHS \\ 0 \leq Y < 1 & RHS \end{cases}$

$$\Rightarrow \underline{E[I(Y \geq 1)]} \leq E[Y]$$

$$\Rightarrow P(Y \geq 1) \leq E[Y] = E\left[\frac{1}{a} |X|\right] = \frac{1}{a} E[|X|].$$

$$\begin{aligned} & \text{①} \\ P\left(\frac{|X|}{a} \geq 1\right) & \leq \\ & \text{②} \\ P(|X| \geq a) & \end{aligned}$$

# Chebyshev's Inequality

$$P(|X-\mu| \geq a) = P(|X-\mu|^2 \geq a^2)$$

Markov's Inequality

$$\leq \frac{1}{a^2} E(|X-\mu|^2) = \frac{1}{a^2} \text{Var}(X)$$

## Theorem

Let  $X$  have mean  $\mu$  and variance  $\sigma^2$ . Then for any  $a > 0$ ,  $P(|X-\mu| \geq a) \leq \frac{\sigma^2}{a^2}$

$$P(|X-\mu| \geq a) \leq \frac{\sigma^2}{a^2}. \quad O\left(\frac{1}{a^2}\right)$$

Application:  $\bar{X}_n$  Sample mean :  $E(\bar{X}_n) = \mu$ ;  $\text{Var}(\bar{X}_n) = \frac{1}{n} \sigma^2$

$$P(|\bar{X}_n - \mu| \geq a) \leq \frac{1}{a^2} \text{Var}(\bar{X}_n) = \frac{\sigma^2}{(n \cdot a^2)} \xrightarrow{n \rightarrow \infty} 0$$

$$\bar{X}_n \xrightarrow{P} \mu.$$

$$O\left(\frac{1}{n}\right)$$

# Proof

## Chernoff's Inequality

$$\forall t > 0 \quad P(X \geq a) = P(e^{tX} \geq e^{ta})$$

Martingale Inequality

$\leq$

$$\frac{E[e^{tX}]}{e^{ta}} f(t)$$

### Theorem

For any r.v.  $X$  and constants  $a > 0$  and  $t > 0$ ,

$$P(X \geq a) \leq \frac{E(e^{tX})}{e^{ta}} f(t) \xrightarrow{\text{MGF.}}$$

$$\forall t > 0 \quad P(X \geq a) \leq f(t)$$

$$\Rightarrow P(X \geq a) \leq \inf_{t > 0} f(t)$$

# Proof

## Chernoff's Technique

$$\forall t < 0 ; P(X \leq a) = P(tx \geq ta)$$

$$= P(e^{tx} \geq e^{ta}) \leq \frac{E[e^{tx}]}{e^{ta}} \cdot f(t)$$

### Theorem

For any r.v.  $X$  and constants  $a$ ,

$$P(X \geq a) \leq \inf_{t>0} \frac{E(e^{tX})}{e^{ta}}$$

$$P(X \leq a) \leq \inf_{t<0} \frac{E(e^{tX})}{e^{ta}}.$$

# Proof

## Example: Normal Distribution

① MGF of  $X$ :  $M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$

②  $P(X > a) \leq \inf_{t > 0} \frac{E[e^{tx}]}{e^{ta}} = \inf_{t > 0} f(t)$

Given  $X \sim \mathcal{N}(\mu, \sigma^2)$ , for arbitrary constant  $a > \mu$ , find the Chernoff bound on  $\underline{P}(X > a)$ .

$$f(t) = \frac{E[e^{tx}]}{e^{ta}} = \frac{M_X(t)}{e^{ta}} = e^{\frac{1}{2}\sigma^2 t^2 + (\mu - \alpha)t}$$
$$= e^{\frac{1}{2}\sigma^2 \left[ \underline{(t + \frac{\mu - \alpha}{\sigma^2})^2} - \frac{(\mu - \alpha)^2}{\sigma^2} \right]} \quad t^* = \frac{\alpha - \mu}{\sigma^2} > 0$$

$$\Rightarrow P(X > a) \leq f(t^*) = e^{-\frac{(\alpha - \mu)^2}{2\sigma^2}}$$

$$\alpha = \mu + \varepsilon$$

$$\mathcal{O}(e^{-\varepsilon^2})$$

$$\Rightarrow P(X > \mu + \varepsilon) \leq e^{-\frac{\varepsilon^2}{2\sigma^2}} \Rightarrow P(X - \mu > \varepsilon) \leq e^{-\frac{\varepsilon^2}{2\sigma^2}}$$

# Solution

# Hoeffding Bound

$$E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \mu$$

## Theorem

Let the random variables  $X_1, X_2, \dots, X_n$  be independent with  $E(X_i) = \mu$ ,  $a \leq X_i \leq b$  for each  $i = 1, \dots, n$ , where  $a, b$  are constants. Then for any  $\epsilon \geq 0$ ,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq \epsilon\right) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}}.$$

$\epsilon \uparrow \downarrow$   
 $n \uparrow \downarrow$

$$\mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) \leq \frac{o(e^{-\epsilon^2})}{\sqrt{n}}$$

# Application: Parameter Estimation

$$\begin{aligned} p \in [\hat{p} - \epsilon, \hat{p} + \epsilon] &\Leftrightarrow \hat{p} - \epsilon \leq p \leq \hat{p} + \epsilon \\ &\Leftrightarrow -\epsilon \leq p - \hat{p} \leq \epsilon \\ &\Leftrightarrow -\epsilon \leq \hat{p} - p \leq \epsilon \Leftrightarrow |\hat{p} - p| \leq \epsilon \end{aligned}$$

Instead of predicting a single value  $\hat{p}$  for the parameter  $p$ , we are given an interval that is likely to contain the parameter:

## Definition

$\delta = 0.05$

A  $1 - \delta$  confidence interval for a parameter  $p$  is an interval  $[\hat{p} - \epsilon, \hat{p} + \epsilon]$  such that

$$Pr(p \in [\hat{p} - \epsilon, \hat{p} + \epsilon]) \geq 1 - \delta.$$

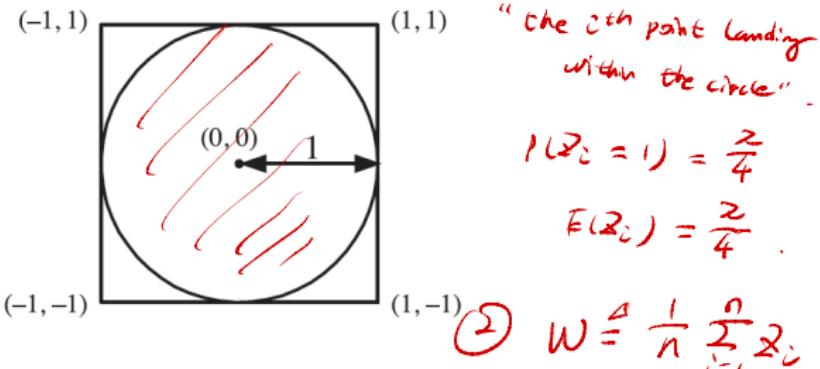
$$\begin{aligned} Pr(|\hat{p} - p| \leq \epsilon) &\geq 1 - \delta \\ \Rightarrow Pr(|\hat{p} - p| > \epsilon) &\leq \underline{\delta} \end{aligned}$$

# Application Example: Monte Carlo Method for Estimation $\pi$

①  $(x_2, y_2)$  ( $-1 \leq x_2 \leq 1, -1 \leq y_2 \leq 1$ )

Circle:  $\{ (x, y) : x^2 + y^2 \leq 1 \}$

$Z_2$ : indicator of the event



②  $W \triangleq \frac{1}{n} \sum_{c=1}^n Z_c$

- A point chosen uniformly at random in the square has probability  $\pi/4$  of landing in the circle

$$E(W) = \frac{\pi}{4}.$$

③  $\hat{\pi} = 4W = 4 \cdot \frac{1}{n} \sum_{c=1}^n Z_c$

Confidence Interval of  $\pi$ .

# Example: Monte Carlo Method for Estimation $\pi$

③  $n \rightarrow \infty$ ,  $\hat{\pi} \rightarrow \pi$  (w.p.1).

$Z_1, \dots, Z_n$  i.i.d.

Param( $\frac{Z}{4}$ )

$$\Pr(|\hat{\pi} - \pi| \geq \varepsilon) = \Pr(|4W - \pi| \geq \varepsilon)$$

$$= \Pr\left(|W - \frac{\pi}{4}| \geq \frac{\varepsilon}{4}\right) = \Pr\left(\left|\frac{1}{n} \sum_{i=1}^n Z_i - E(Z)\right| \geq \frac{\varepsilon}{4}\right)$$

~~Hoeffding's inequality~~

$$\frac{\text{Var}\left(\frac{1}{n} \sum_{i=1}^n Z_i\right)}{\left(\frac{\varepsilon}{4}\right)^2} = \frac{\frac{2}{n}(1-\frac{2}{n})}{\left(\frac{\varepsilon}{4}\right)^2 \cdot n}$$

Hoeffding's inequality

$$\text{bound} \leq 2 \cdot e^{-\frac{2n(\frac{\varepsilon}{4})^2}{(1-\frac{2}{n})^2}} = \frac{2e^{-\frac{1}{8}n\varepsilon^2}}{n} = \delta$$

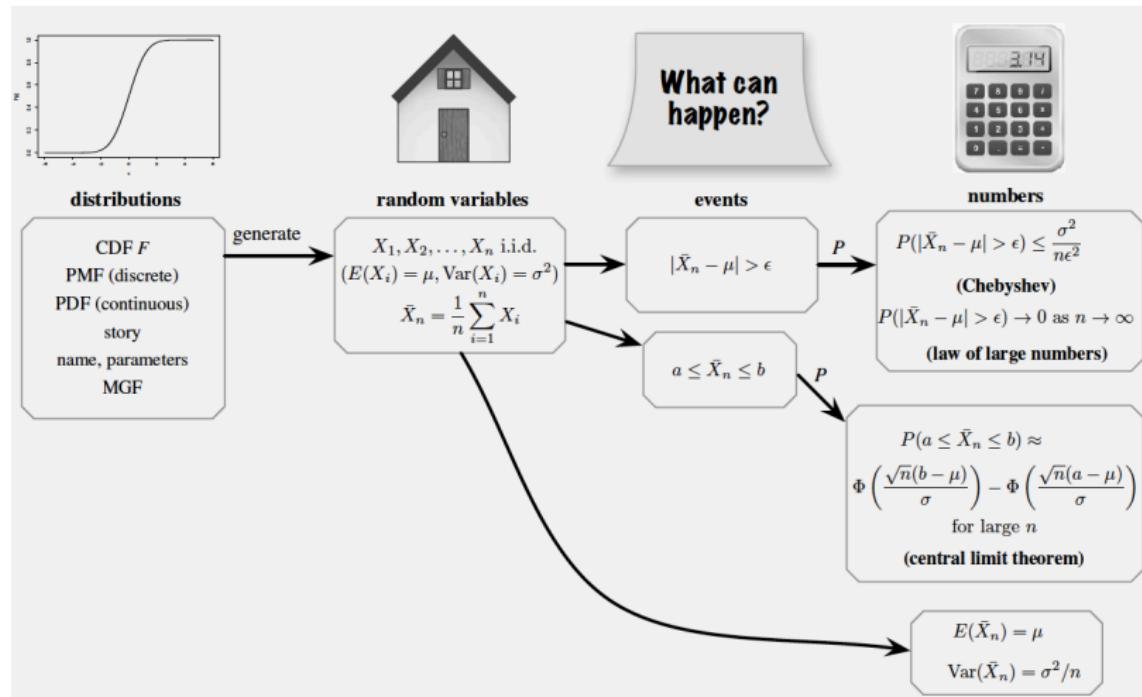
$$\Rightarrow \varepsilon = \sqrt{\frac{8 \ln(\frac{1}{\delta})}{n}}$$

$\delta = 0.05$

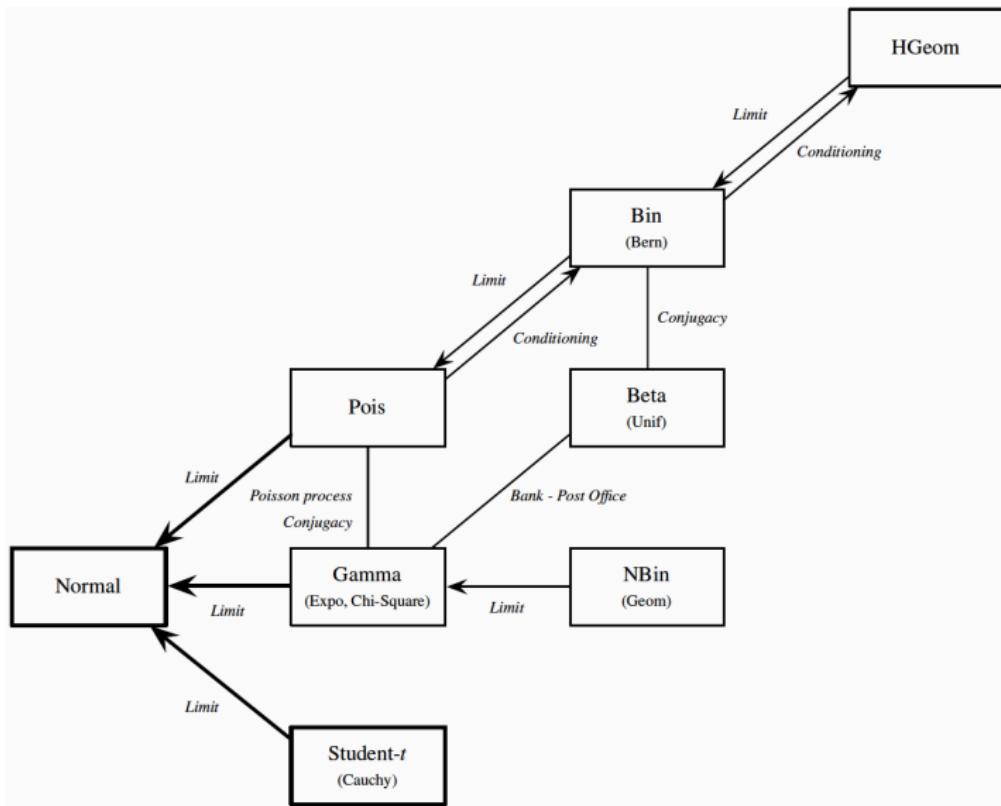
$$\Rightarrow \Pr\left(\pi \in \left(\hat{\pi} - \sqrt{\frac{8 \ln(\frac{1}{\delta})}{n}}, \hat{\pi} + \sqrt{\frac{8 \ln(\frac{1}{\delta})}{n}}\right)\right) \geq 1 - \delta$$

# Example: Monte Carlo Method for Estimation $\pi$

# Summary 1



# Summary 2



# References

- Chapter 10 of **BH**
- Chapter 5 of **BT**