

# Lecture 1: Probability & Statistics Review

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

February 26, 2025

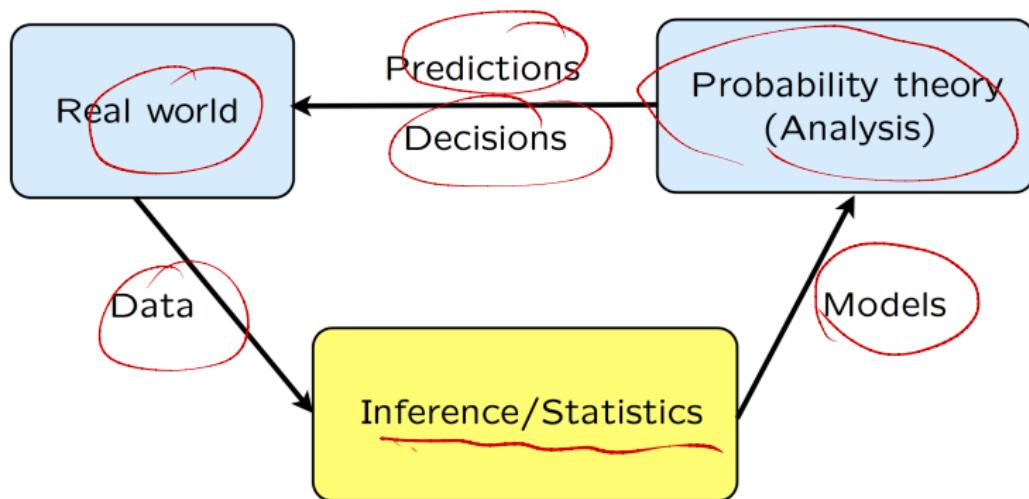
# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

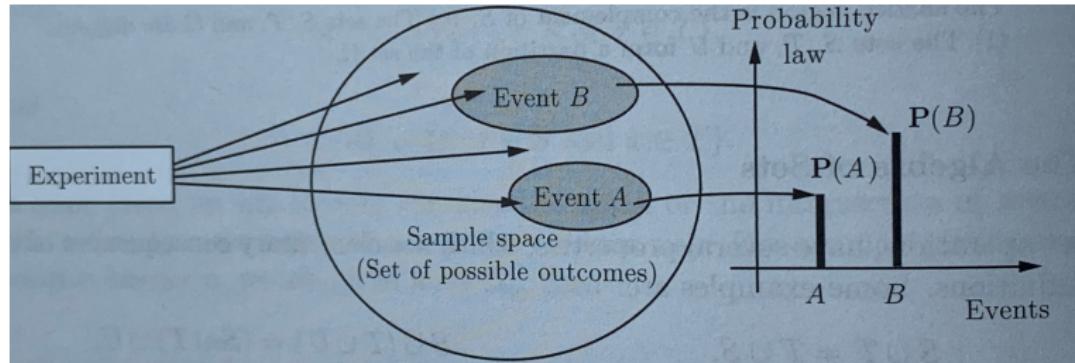
# Outline

- 1 **Stochastic Laws in Random Worlds**
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# Understanding Random Worlds: Framework



# Understanding Random Worlds: Modeling



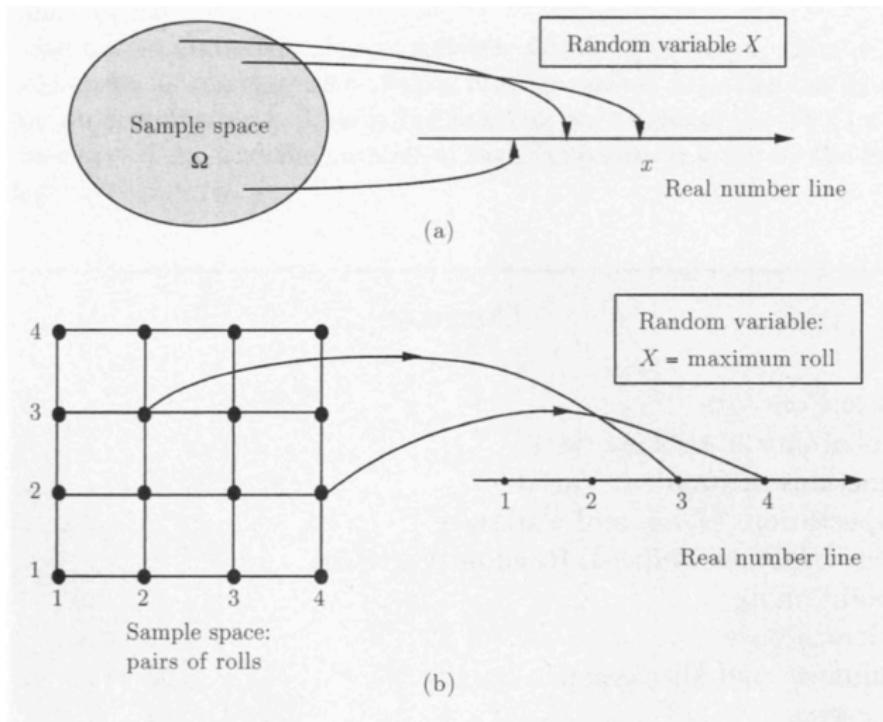
# Bayes' Rule

## Theorem

*For any events  $A$  and  $B$  with positive probabilities,*

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}.$$

# Understanding Random Worlds: Modeling



# General Bayes' Rule

|                | $Y$ discrete                                       | $Y$ continuous                                     |
|----------------|--|--|
| $X$ discrete   | $P(Y = y X = x) = \frac{P(X=x Y=y)P(Y=y)}{P(X=x)}$ | $f_Y(y X = x) = \frac{P(X=x Y=y)f_Y(y)}{P(X=x)}$   |
| $X$ continuous | $P(Y = y X = x) = \frac{f_X(x Y=y)P(Y=y)}{f_X(x)}$ | $f_{Y X}(y x) = \frac{f_{X Y}(x y)f_Y(y)}{f_X(x)}$ |

# General LOTP

|                | $Y$ discrete                               | $Y$ continuous  |
|----------------|--|---|
| $X$ discrete   | $P(X = x) = \sum_y P(X = x Y = y)P(Y = y)$ | $P(X = x) = \int_{-\infty}^{\infty} P(X = x Y = y)f_Y(y)dy$ |
| $X$ continuous | $f_X(x) = \sum_y f_X(x Y = y)P(Y = y)$     | $f_X(x) = \int_{-\infty}^{\infty} f_{X Y}(x y)f_Y(y)dy$     |

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# Important Concepts

- Population: distribution  $\mathcal{F}$
- Sample: random vector  $\mathbf{X} = (X_1, \dots, X_n)$  where  $n$  is the sample size
- Random Sample:  $\{X_i\}$  are i.i.d random variables and  $X_i \sim \mathcal{F}$
- Data: real vector  $\mathbf{x} = (x_1, \dots, x_n)$ , the value of sample  $\mathbf{X}$
- From sample to infer property of Population
- Statistic: a function of sample  $\mathbf{X}$

$$\bar{X} = \frac{1}{n}(x_1 + \dots + x_n)$$

# Important Concepts

Difference between an estimate and an estimator

- An estimator is a function of the sample
- An estimator is a function of the random variables  $X_1, \dots, X_n$ .
- An estimate is the realized value of an estimator (a number) that is obtained when a sample is actually taken
- An estimate is a function of the realized values  $x_1, \dots, x_n$ .

# Nonparametric Models



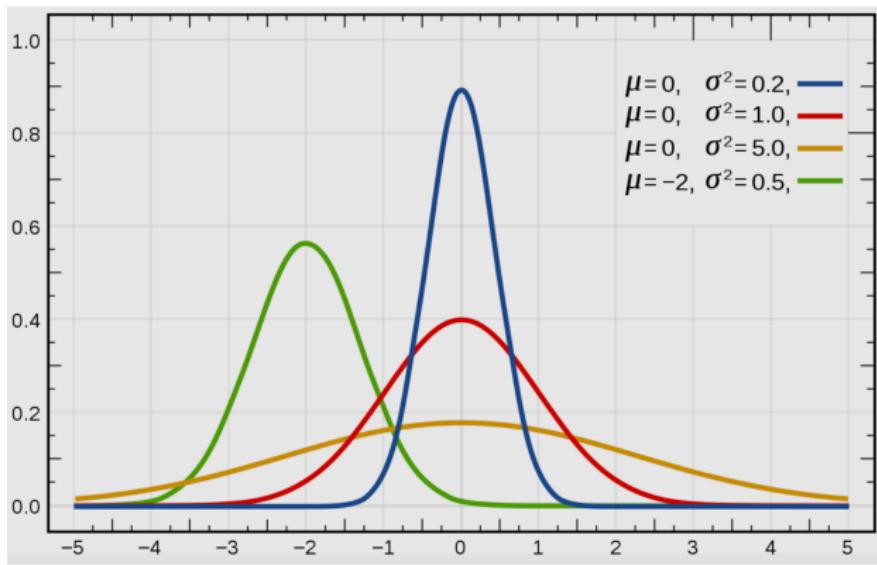
- A nonparametric model is a statistical model that cannot be parameterized by a finite number of parameters.
- Example: (nonparametric estimation of the CDF) Let  $X_1, \dots, X_n$  be independent observations from a CDF  $F$ . Our task is to estimate  $F$ .

# Our Focus: Parameterized Statistical Inference

- PMT or PDF*
- Given a parametric distribution model (a family of PMFs or PDFs)  $\mathcal{F} = \{\rho(x; \theta) : \theta \in \mathcal{R}\}$
  - $\theta$  is an unknown parameter in a parameter space  $\mathcal{R}$
  - Now given (random)sample from such model:  $\mathbf{X} = (X_1, \dots, X_n)$
  - How to make parameterized statistical inference?

# Example: Parameterized Distribution Model

$$\left\{ p(x; \underline{\mu}, \underline{\sigma}) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right), \mu \in \mathbb{R}, \sigma > 0 \right\}.$$

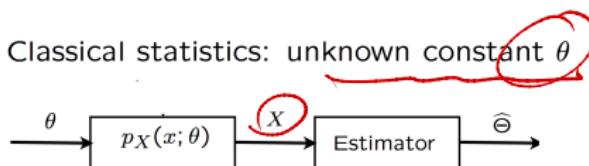


# Parameterized Statistical Inference: Bayesian versus Frequentist

- Difference relates to the nature of the unknown parameter  $\theta$
- Treated as a random variable  $\Theta$  with prior (known) distribution:  
**Bayesian approach**
- Treated as an unknown constant  $\theta$ : **frequentist approach**

# Statistical Inference: Frequentist Perspective

- Classical statistics: unknown constant  $\theta$



- Hypothesis testing:  $H_0 : \theta = 1/2$  versus  $H_1 : \theta = 3/4$
- Composite hypotheses:  $H_0 : \theta = 1/2$  versus  $H_1 : \theta \neq 1/2$
- Estimation: design an **estimator**  $\hat{\theta}$ , to “keep estimation **error**  $\hat{\theta} - \theta$  small”

# Estimation Method: Maximum Likelihood Estimation (MLE)

- We observe a particular data  $\mathbf{x} = (x_1, \dots, x_n)$ ,
- Likelihood: the probability (or probability density) of seeing data  $\mathbf{x}$  under different values of parameter  $\theta$ , i.e.,  $p(\mathbf{x}; \theta)$
- A **maximum likelihood estimate** (MLE) is a value of the parameter  $\theta$  that maximizes the likelihood  $p(\mathbf{x}; \theta)$  over all possible values:

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{x}; \theta)$$

# MLE under Independent Case

- Random Sample:  $\{X_i\}$  are i.i.d., we have

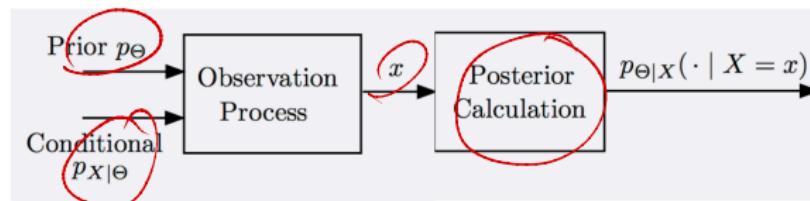
$$\log[p(\mathbf{x}; \theta)] = \log \prod_{i=1}^n p(x_i; \theta) = \sum_{i=1}^n \log[p(x_i; \theta)]$$

- Thus a **maximum likelihood estimate** (MLE) under independent case is shown as follows:

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} p(\mathbf{x}; \theta) = \arg \max_{\theta} \log[p(\mathbf{x}; \theta)] \\ &= \arg \max_{\theta} \sum_{i=1}^n \log[p(x_i; \theta)]\end{aligned}$$

# Statistical Inference: Bayesian Perspective

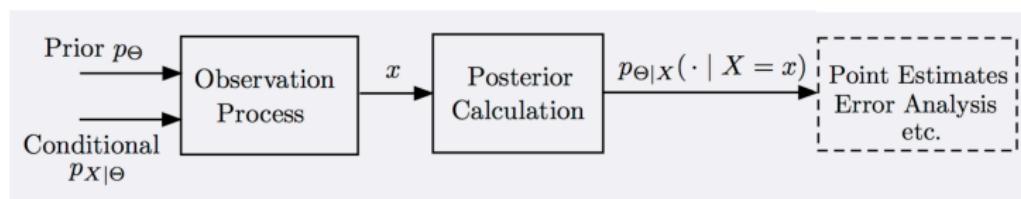
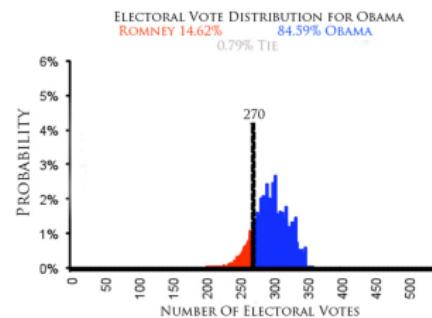
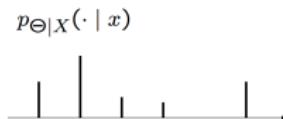
- Unknown  $\Theta$ 
    - treated as a random variable
    - prior distribution  $p_\Theta$  or  $f_\Theta$
  - Observation  $X$ 
    - observation model  $p_{X|\Theta}$  or  $f_{X|\Theta}$
  - Use appropriate version of the Bayes rule to find  $p_{\Theta|X}(\cdot | X = x)$  or  $f_{\Theta|X}(\cdot | X = x)$
- Literally*



# The Output of Bayesian Statistical Inference

The complete answer is a posterior distribution:

PMF  $p_{\Theta|X}(\cdot | x)$  or PDF  $f_{\Theta|X}(\cdot | x)$



# General Bayes' Rule for Bayesian Statistical Inference

## The Four Versions of Bayes' Rule

- $\Theta$  discrete,  $X$  discrete:

$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta)p_{X|\Theta}(x | \theta)}{\sum_{\theta'} p_\Theta(\theta')p_{X|\Theta}(x | \theta')}.$$

$\propto p_\Theta(\theta) \cdot p_{X|\Theta}(x | \theta)$

- $\Theta$  discrete,  $X$  continuous:

$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta)f_{X|\Theta}(x | \theta)}{\sum_{\theta'} p_\Theta(\theta')f_{X|\Theta}(x | \theta')}.$$

- $\Theta$  continuous,  $X$  discrete:

$$f_{\Theta|X}(\theta | x) = \frac{f_\Theta(\theta)p_{X|\Theta}(x | \theta)}{\int f_\Theta(\theta')p_{X|\Theta}(x | \theta') d\theta'}.$$

- $\Theta$  continuous,  $X$  continuous:

$$f_{\Theta|X}(\theta | x) = \frac{f_\Theta(\theta)f_{X|\Theta}(x | \theta)}{\int f_\Theta(\theta')f_{X|\Theta}(x | \theta') d\theta'}.$$

# Bayes' Rule: Bayesian Perspective

## Theorem

*A represents cognition, B represents data, then*

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}.$$

- $P(A)$ : Prior Cognition
- $P(B)$ : Evidence
- $P(B|A)$ : Likelihood (info of data given prior cognition)
- $P(A|B)$ : Posterior Cognition

# Estimation Method I

- The Maximum A Posteriori Probability (MAP)
- Given the observation value  $x$ , the MAP rule selects a value  $\hat{\theta}$  that maximizes the posterior probability (probability density)  $p_{\Theta|x}(\theta|x)$ :

$$\hat{\theta} = \arg \max_{\theta} p_{\Theta|x}(\theta|x)$$

- Equivalently,

$$\hat{\theta} = \arg \max_{\theta} p_{\Theta}(\theta) p_{X|\Theta}(x|\theta)$$

# Estimation Method II

- Posterior Mean: Given the observation data  $\mathbf{x}$ , the estimate of  $\Theta$  is  $\hat{\theta} = \underline{E[\Theta | \mathbf{X} = \mathbf{x}]}$
- Now given the observation sample  $\mathbf{X}$ , the estimator of  $\Theta$  is  $\hat{\theta} = \underline{E[\Theta | \mathbf{X}]}$

## Outline

- 1 Stochastic Laws in Random Worlds
  - 2 Statistical Inference
  - 3 Conjugate Prior: A Weapon of Bayesian  
A Weapon of Bayesian
  - 4 Conditional Expectation
  - 5 Monte Carlo Method: Another Weapon of Bayesian
  - 6 Sampling: Random Variable Generation
  - 7 Sampling: Random Vector Generation
  - 8 Monte Carlo Integration
  - 9 Performance Analysis of Monte Carlo Integration
  - 10 References

# Conjugate Prior

- Before Monte Carlo, posterior calculation is hard
- Conjugate Prior: reduce the computing complexity of posterior distribution
- Loosely speaking, a prior distribution is conjugate to the likelihood model if both the prior and posterior distribution stay in the same distribution family.

## Story: Beta-Binomial Conjugacy

We have a coin that lands Heads with probability  $p$ , but we don't know what  $p$  is. Our goal is to infer the value of  $p$  after observing the outcomes of  $n$  tosses of the coin. The larger that  $n$  is, the more accurately we should be able to estimate  $p$ .

# Bayesian Inference

- Treats all unknown quantities as random variables.
- In the Bayesian approach, we would treat the unknown probability  $p$  as a random variable and give  $p$  a distribution.
- This is called a prior distribution, and it reflects our uncertainty about the true value of  $p$  before observing the coin tosses.
- After the experiment is performed and the data are gathered, the prior distribution is updated using Bayes' rule; this yields the posterior distribution, which reflects our new beliefs about  $p$ .
- Now we adopt the beta distribution as the prior distribution.

# Beta Distribution

$$\frac{a=1, b=1;}{\text{Unif}(0,1)}$$

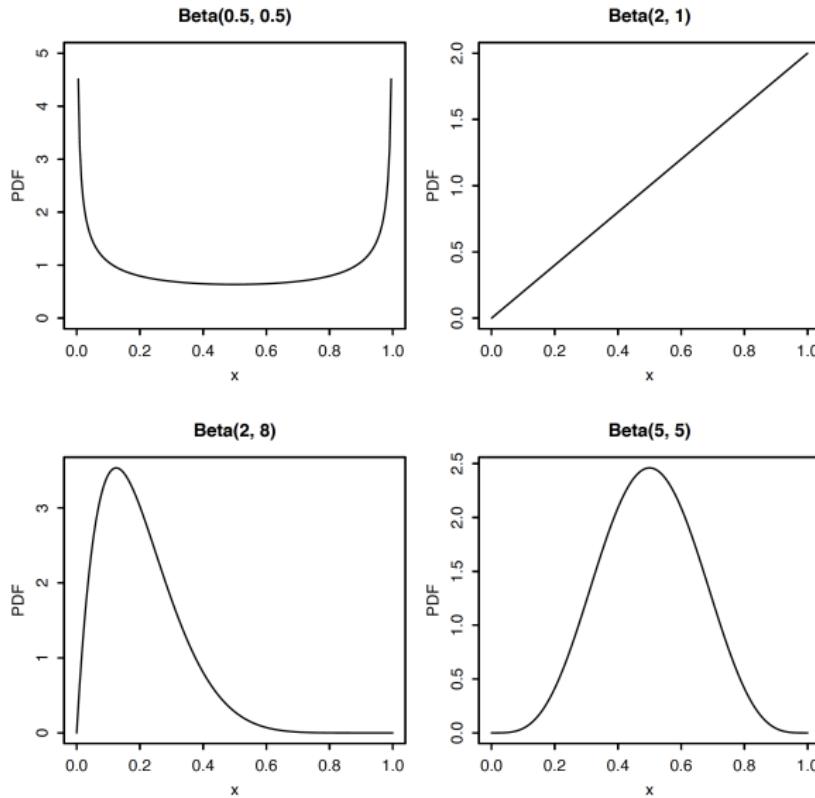
## Definition

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. It is denoted as  $X \sim \text{Beta}(a, b)$ . Beta distribution is a generalization of uniform distribution.

# PDF of Beta Distribution



# Story: Beta-Binomial Conjugacy

$$\textcircled{2} \quad f_p(p | X=k) = \frac{\Pr(X=k | P=p)}{\Pr(X=k)} \cdot \underline{f_p(p)}$$

$$= \frac{\binom{n}{k} p^k (1-p)^{n-k}}{\Pr(X=k)} \cdot \frac{1}{\text{Bayes}} \cdot \underline{p^{a-1} (1-p)^{b-1}}$$

\textcircled{3}  $f_p(p | X=k)$ : a function of  $p$

$$f_p(p | X=k) \propto \underline{p^k (1-p)^{n-k} \cdot p^{a-1} (1-p)^{b-1}}$$

$$\sim \text{Beta}(a+b, n-k+b)$$

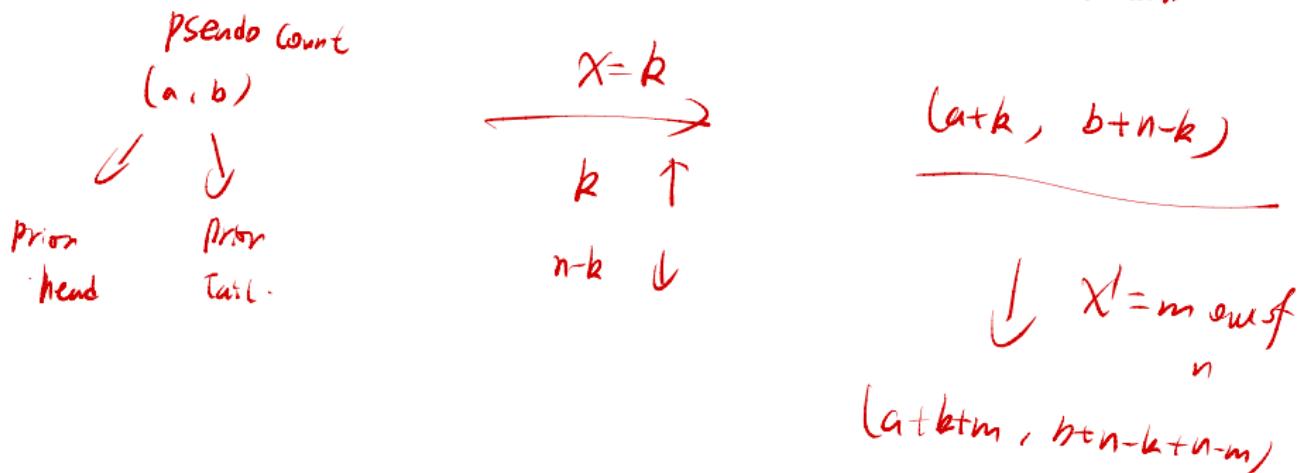
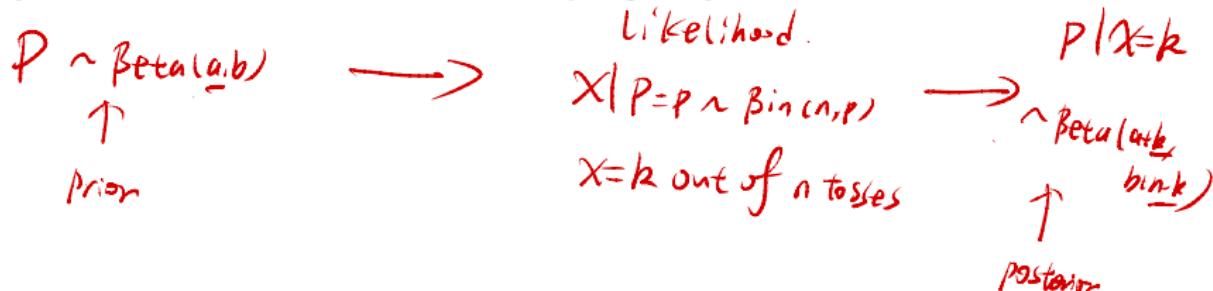
①  $P = r.v.$  Prior Beta  
 $X$ : # of heads out of  $n$  tosses  
 $X | P=p \sim \text{Bin}(n, p)$ .  
Likelihood

$$\Pr(X=k) = \int_0^1 \Pr(X=k | P=p) f_p(p) dp$$

$$= p^{k+a-1} (1-p)^{n-k+b-1}$$

$0 < p < 1$

# Story: Beta-Binomial Conjugacy



# Story: Beta-Binomial Conjugacy

# Story: Beta-Binomial Conjugacy

- Furthermore, notice the very simple formula for updating the distribution of  $p$ .
- We just add the number of observed successes,  $k$ , to the first parameter of the Beta distribution.
- We also add the number of observed failures,  $n - k$ , to the second parameter of the Beta distribution.
- So  $a$  and  $b$  have a concrete interpretation in this context:
  - ▶  $a$  as the number of prior successes in earlier experiments
  - ▶  $b$  as the number of prior failures in earlier experiments
  - ▶  $a, b$ : pseudo counts

# Mean vs. Bayesian(Posterior) Average

- Infer the value of  $p$  (probability of coin lands heads)
- Observed  $k$  heads out of  $n$  tosses of the coin
- Mean:  $\frac{k}{n}$  *MLE*
- Bayesian(Posterior) Average:  $E(p|X = k) = \frac{a+k}{a+b+n}$
- Suppose the prior distribution is  $\text{Unif}(0,1)$ :  $a = 1, b = 1$
- Bayesian(Posterior) Average:  $\frac{k+1}{n+2}$
- When  $k = n$ , we have:  $1$  (mean) vs.  $\frac{n+1}{n+2}$  (Bayesian average)

MLE

$n \rightarrow \infty$

$n=3$

## Story: Beta-Binomial Conjugacy

If we have a Beta prior distribution on  $p$  and data that are conditionally Binomial given  $p$ , then when going from prior to posterior, we don't leave the family of Beta distributions. We say that **the Beta distribution is the conjugate prior of the Binomial distribution.**

# Likelihood Model: Discrete

| Sample Space                    | Sampling Dist.          | Conjugate Prior          | Posterior   |
|---------------------------------|-------------------------|--------------------------|---|
| $\mathcal{X} = \{0, 1\}$        | Bernoulli( $\theta$ )   | Beta( $\alpha, \beta$ )  | Beta( $\alpha + n\bar{X}, \beta + n(1 - \bar{X})$ ) |
| $\mathcal{X} = \mathbb{Z}_+$    | Poisson( $\lambda$ )    | Gamma( $\alpha, \beta$ ) | Gamma( $\alpha + n\bar{X}, \beta + n$ )             |
| $\mathcal{X} = \mathbb{Z}_{++}$ | Geometric( $\theta$ )   | Gamma( $\alpha, \beta$ ) | Gamma( $\alpha + n, \beta + n\bar{X}$ )             |
| $\mathcal{X} = \mathbb{H}_K$    | Multinomial( $\theta$ ) | Dirichlet( $\alpha$ )    | Dirichlet( $\alpha + n\bar{X}$ )                    |

# Likelihood Model: Continuous

| Sampling Dist.                        | Conjugate Prior                        | Posterior   |
|---------------------------------------|--|---|
| Uniform( $\theta$ )                   | Pareto( $\nu_0, k$ )                   | Pareto $(\max\{\nu_0, X_{(n)}\}, n + k)$  |
| Exponential( $\theta$ )               | Gamma( $\alpha, \beta$ )               | Gamma( $\alpha + n, \beta + n\bar{X}$ )   |
| $N(\mu, \sigma^2)$ , known $\sigma^2$ | $N(\mu_0, \sigma_0^2)$                 | $N \left( \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1} \left( \frac{\mu_0}{\sigma_0^2} + \frac{n\bar{X}}{\sigma^2} \right), \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1} \right)$ |
| $N(\mu, \sigma^2)$ , known $\mu$      | InvGamma( $\alpha, \beta$ )            | InvGamma $\left( \alpha + \frac{n}{2}, \beta + \frac{n(X - \mu)^2}{2} \right)$  |
| $N(\mu, \sigma^2)$ , known $\mu$      | ScaledInv- $\chi^2(\nu_0, \sigma_0^2)$ | ScaledInv- $\chi^2 \left( \nu_0 + n, \frac{\nu_0 \sigma_0^2}{\nu_0 + n} + \frac{n(\bar{X} - \mu)^2}{\nu_0 + n} \right)$   |
| $N(\mu, \Sigma)$ , known $\Sigma$     | $N(\mu_0, \Sigma_0)$                   | $N \left( \mathbf{K} \left( \Sigma_0^{-1} \mu_0 + n \Sigma^{-1} \bar{X} \right), \mathbf{K} \right)$ , $\mathbf{K} = (\Sigma_0^{-1} + n \Sigma^{-1})^{-1}$  |
| $N(\mu, \Sigma)$ , known $\mu$        | InvWishart( $\nu_0, \mathbf{S}_0$ )    | InvWishart( $\nu_0 + n, \mathbf{S}_0 + n\bar{\mathbf{S}}$ ), $\bar{\mathbf{S}}$ sample covariance   |

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# Motivation

- Conditional expectation is a powerful tool for calculating expectations: first-step analysis
- Conditional expectation allows us to predict or estimate unknowns based on whatever evidence is currently available.
- Conditional Expectation given an event:  $E(Y|A)$
- Conditional Expectation given a random variable:  $E(Y|X)$

# Conditional Expectation Given An Event

## Definition

Let  $A$  be an event with positive probability. If  $Y$  is a discrete r.v., then the *conditional expectation of  $Y$  given  $A$*  is

$$E(Y|A) = \sum_y yP(Y=y|A),$$

where the sum is over the support of  $Y$ . If  $Y$  is a continuous r.v. with PDF  $f$ , then

$$E(Y|A) = \int_{-\infty}^{\infty} yf(y|A) dy,$$

where the conditional PDF  $f(y|A)$  is defined as the derivative of the conditional CDF  $F(y|A) = P(Y \leq y|A)$ , and can also be computed by a hybrid version of Bayes' rule:

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

# Intuition for $E(Y|A)$

$$E(Y) \approx \frac{1}{n} (Y_1 + \dots + Y_n)$$

## Principle

$E(Y|A)$  is approximately the average of  $Y$  in a large number of simulation runs in which  $A$  occurred.

# Life Expectancy

$$E[T] = 70;$$

$$E[T \mid T \geq 20]$$

# Law of Total Expectation

## Theorem

Let  $A_1, \dots, A_n$  be a partition of a sample space, with  $P(A_i) > 0$  for all  $i$ , and let  $Y$  be a random variable on this sample space. Then

$$\underline{E(Y)} = \sum_{i=1}^n \underline{E(Y|A_i) P(A_i)}.$$

# Conditional Expectation Given An R.V.

## Definition

Let  $g(x) = E(Y|X = x)$ . Then the *conditional expectation of Y given X*, denoted  $E(Y|X)$ , is defined to be the random variable  $g(X)$ . In other words, if after doing the experiment  $X$  crystallizes into  $x$ , then  $E(Y|X)$  crystallizes into  $g(x)$ .

## Remark

- $E(Y|X)$  is a function of  $X$ , and it is a random variable.
- It makes sense to compute  $E(E(Y|X))$  and  $\text{Var}(E(Y|X))$ .

# Property I: Dropping What's Independent

$$g(x) = E[Y|X=x] = \underline{E[Y]}, \forall x.$$
$$\Rightarrow g(x) = \underline{E[Y]}$$

$$\Rightarrow E[Y|X] = g(x) = \underline{E[Y]}$$

## Theorem

If  $X$  and  $Y$  are independent, then  $E(Y|X) = E(Y)$ .

## Property II: Taking Out What's Known

$$1^o. \quad g(x) = E[h(x)Y | X=x]$$

$$= E[\underline{h(x)}Y | X=x] = \underline{h(x)} E[\underline{Y | X=x}]$$

$$2^o. \quad g(X) = h(X) \cdot E[Y|X]$$

### Theorem

For any function  $h$ ,

$$E(h(X) Y | X) = \underline{h(X)} E(Y | X)$$

## Property III: Linearity

$$\begin{aligned}g(x) &= E[Y_1 + Y_2 | X=x] = E[Y_1 | X=x] + E[Y_2 | X=x] \\ \Rightarrow g(x) &= E(Y_1 | x) + E(Y_2 | x)\end{aligned}$$

### Theorem

$$\underline{E(Y_1 + Y_2 | X)} = \underline{E(Y_1 | X)} + \underline{E(Y_2 | X)}.$$

## Adam's Law

W.L.O.G.  $X$  and  $Y$  are both discrete r.v.s.

$$1^{\circ} \quad g(X) = E[Y|X] :$$

$$\underline{g(X)} = E[Y|X=x] = \sum_y y \cdot P(Y=y|X=x)$$

$$2^{\circ}. \quad LHS : E[E[Y|X]] = E[g(X)] \\ = \sum_x g(x) \cdot P(X=x)$$

## Theorem (Law of Iterated Expectation)

For any r.v.s  $X$  and  $Y$ ,

$$\begin{array}{c} \text{LHS} \\ E(\underline{E(Y|X)}) = \underline{E(Y)}. \end{array}$$

$$= \sum_x \left[ \sum_y y \cdot P(Y=y|X=x) \right] \cdot P(X=x)$$

$$= \sum_y y \cdot \left[ \sum_x P(Y=y|X=x) \cdot P(X=x) \right] = \sum_y y \cdot \left[ \sum_x P(Y=y, X=x) \right]$$

$$= \sum_y y \cdot \underline{P(Y=y)} = \underline{E[Y]}$$

# Proof

## Adam's Law with Extra Conditioning

$$\hat{E} = E(\cdot | z)$$

Adam's Law :  $E[E(Y|X)] = E(Y)$

$$\hat{E}[\hat{E}(Y|X)] = \hat{E}(Y)$$

$$\hat{E}(Y) = E(Y|z)$$

Theorem (Law of Iterated Expectation with Extra Conditioning)

For any r.v.s  $X, Y, Z$ , we have

$$\begin{aligned} & \hat{E}(\hat{E}(Y|X)) \\ &= \hat{E}(E(Y|X, Z)) \\ &= E[\hat{E}(Y|X, Z)|Z] \\ & E(E(Y|X, Z)|Z) = E(Y|Z) \\ & E(E(X|Z, Y)|Y) = E(X|Y) \end{aligned}$$

$$\hat{E} = E(\cdot | Y)$$

# Conditional Variance

$$\text{Var}(Y) = E[(Y - \underline{E}(Y))^2]$$

$$\widehat{E}(\cdot) = E(\cdot | X)$$

## Definition

$$\text{Var}(Y|X) = \widehat{E}[(Y - \widehat{E}(Y))^2]$$

The *conditional variance of Y given X* is

$$\text{Var}(Y|X) = \underline{E}\left((Y - E(Y|X))^2 | X\right).$$

This is equivalent to

$$\text{var}(Y) = \underline{E}(Y^2) - \underline{E}^2(Y)$$

$$\text{Var}(Y|X) = \underline{E}(Y^2|X) - (\underline{E}(Y|X))^2.$$

$$\text{var}(Y|X) = \widehat{E}(Y^2) - \widehat{E}^2(Y)$$

# Eve's law

## Theorem (Law of Total Variance)

For any r.v.s  $X$  and  $Y$ ,

$$\text{Var}(Y) = \underbrace{E(\text{Var}(Y|X))}_{\text{EV}} + \underbrace{\text{Var}(E(Y|X))}_{\text{VE}}.$$

*The ordering of E's and Var's on the right-hand side spells EVVE, whence the name Eve's law. Eve's law is also known as the law of total variance or the variance decomposition formula.*

$$\text{Proof } \textcircled{1} \ g(x) = E[Y|x] \Rightarrow \underline{E[g(x)]} = E[\underline{E[Y|x]}] = \underline{E[Y]}$$

$$\begin{aligned} \textcircled{2} \ \underbrace{E[\text{Var}(Y|x)]} &= E[E[Y^2|x] - \underline{(E(Y|x))^2}] \\ &= E[E[Y^2|x]] - g^2(x) = \underline{E[E[Y^2|x]]} - \underline{E[g^2(x)]} \\ &= \underline{E[Y^2]} - \cancel{\underline{E[g^2(x)]}} \end{aligned}$$

$$\begin{aligned} \textcircled{3} \ \underbrace{\text{Var}[E(Y|x)]} &= \text{Var}[g(x)] = E[g^2(x)] - \underline{E^2[g(x)]} \\ &= \cancel{E[g^2(x)]} - \cancel{E^2(Y)} \quad \textcircled{④} \end{aligned}$$

$$\begin{aligned} \textcircled{2} + \textcircled{3} \ E[\text{Var}(Y|x)] + \text{Var}[E(Y|x)] &= E[Y^2] - E^2(Y) \\ &= \text{Var}(Y) \end{aligned}$$

# Prediction & Estimation

- Estimate  $Y$  from the observed value  $X$
- Choose the estimator (inference function)  $\underline{g(\cdot)}$  to minimize the expected error  $E(c(Y, g(X)))$
- $c(Y, \hat{Y})$  is the cost of guessing  $\hat{Y}$  when the actually value is  $\underline{Y}$ .
- When  $c(Y, \hat{Y}) = ||Y - \hat{Y}||^2$ , the best guess is called “the least square estimate (LSE)” estimate of  $Y$  given  $X$ .
- Further, if the function  $\underline{g(\cdot)}$  is restricted to be linear, i.e., of the form  $a + bX$ , it is called “the Linear Least Square Estimate (LLSE)” estimate of  $Y$  given  $X$ .
- Further, if the function  $\underline{g(\cdot)}$  can be arbitrary, it is called “the Minimum Mean Square Estimate (MMSE)” estimate of  $Y$  given  $X$ .

# Linear Least Square Estimate

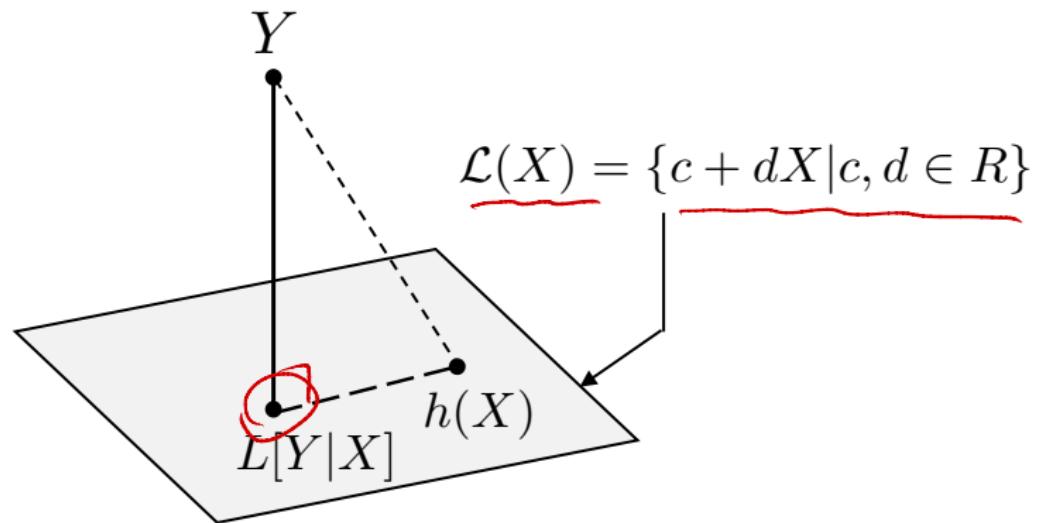
$$\min_{a,b} \underline{E[(Y - a - bX)^2]}$$

## Theorem

The Linear Least Square Estimate (LLSE) of  $Y$  given  $X$ , denoted by  $L[Y|X]$ , is the linear function  $\underline{a + bX}$  that minimizes  $E[(Y - a - bX)^2]$ .  
In fact,

$$L[Y|X] = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)}(X - E(X))$$

# LLSE: Projection Perspective



# Minimum Mean Square Error Estimator

$$\text{LLSE} \quad \min_{\hat{Y}} E[(Y - \hat{Y})^2], \quad \hat{Y} = a + bX \\ \Rightarrow \hat{Y}^* = L(Y|X)$$

$$\text{MMSE} \quad \min_{\hat{Y}} E[(Y - \hat{Y})^2], \quad \hat{Y} = g(X)$$

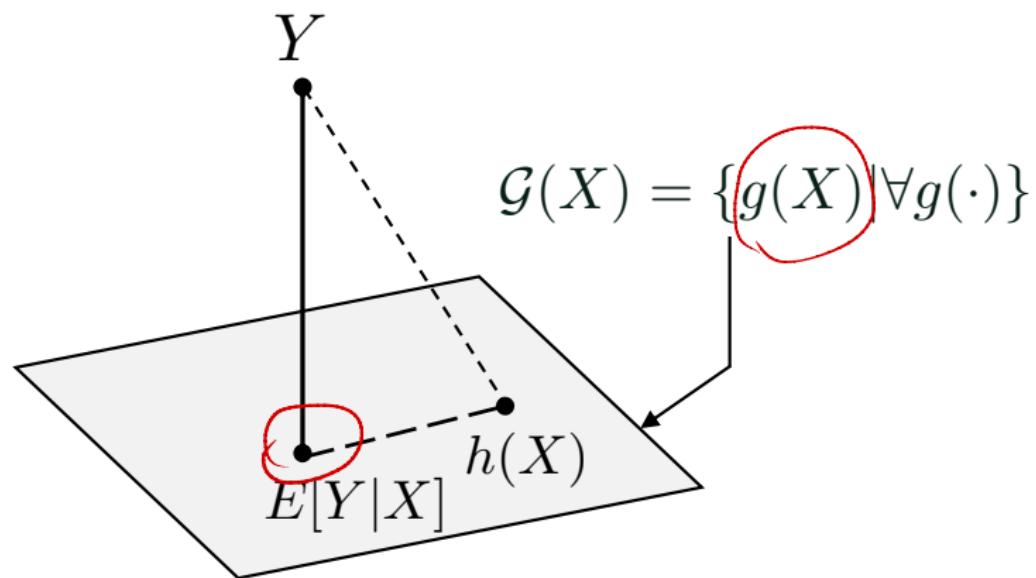
Theorem

The MMSE of  $Y$  given  $X$  is given by

$$\Rightarrow \hat{Y}^* = E[Y|X]$$

$$g(X) = E[Y|X]$$

# MMSE: Projection Perspective



# Projection Interpretation

$$\text{dist}(x, Y) = \sqrt{x \cdot Y}$$
$$= \sqrt{E(x \cdot Y)}$$

$E[X^2] < \infty$  ① Inner product

$$\cos \theta = \frac{\langle x, Y \rangle}{\|x\| \cdot \|Y\|} \quad \langle x, Y \rangle = E[x \cdot Y]$$
$$\|x\| = \sqrt{E[x^2]}$$
$$= \sqrt{E[X^2]}$$

$x \perp Y$

Theorem ②  $x$  and  $Y$  are orthogonal if  $\langle x, Y \rangle = 0 \Leftrightarrow E[x \cdot Y] = 0$

For any function  $h$ , the r.v.  $Y - E(Y|X)$  is uncorrelated with  $h(X)$ .

Equivalently,

$$E((Y - E(Y|X))h(X)) = 0.$$

$(Y - E(Y|X)) \perp h(X)$

(This is equivalent since  $E(Y - E(Y|X)) = 0$ , by linearity and Adam's law.)

$$= E(Y) - E(E(Y|X)) = E(Y) - E(Y) = 0$$

③  $\text{Cov}(x, Y) = E(x \cdot Y) - E(x) \cdot E(Y)$

If  $E(x) = 0$  or  $E(Y) = 0$  or both,  $\text{Cov}(x, Y) = E(x \cdot Y)$   
 $\Rightarrow$  Uncorrelated  $\Leftrightarrow$  orthogonal

## Proof

$$E[(Y - E[Y|X]) \cdot h(x)]$$

$$= E[Yh(x) - E[Y|X] \cdot h(x)]$$

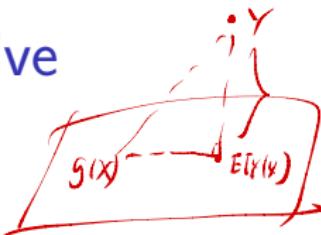
$$= \underline{E[Yh(x)]} - \underline{E[E[Y|X] \cdot h(x)]}$$

$$\quad \quad \quad - \underline{\underline{E[E[Yh(x)|X]]}}$$

$$\quad \quad \quad - \underline{\underline{E[Yh(x)]}}$$

$$= 0$$

# Prediction Perspective


$$\begin{aligned} & E[(Y - E[Y|X])^2] \\ & \leq \underline{\overline{E(Y - g(X))^2}} \\ & \quad \vee g(\cdot) \end{aligned}$$

- Predict or estimate the future observations or unknown parameters based on data
- $E(Y|X)$  is our best predictor of  $Y$  based on  $X$ .
- Best means it is the function of  $X$  with the lowest mean squared error (expected squared difference between  $Y$  and prediction of  $Y$ ).

Proof 1<sup>o</sup>.  $\hat{Y}$ : estimator of  $Y$  based on  $X$ . ( $\hat{Y} = g(X)$ )

$$E[(Y - \hat{Y})^2] = E[(Y - g(X))^2], \because Y - g(X) = \underbrace{Y - E[Y|X]}_A + \underbrace{E[Y|X] - g(X)}_B$$

$$2^o. \underbrace{E[(Y - g(X))^2]}_{E[AB]} = E[(A+B)^2] = \underbrace{E[A^2]}_{E[A^2]} + \underbrace{2E[AB]}_{2E[AB]} + \underbrace{E[B^2]}_{E[B^2]}$$

$$E[AB] = E[(Y - E[Y|X])(\underbrace{E[Y|X] - g(X)}_{\phi(X)})] = 0$$

$$= \underbrace{E[(Y - E[Y|X])^2]}_{=} + \underbrace{E[(E[Y|X] - g(X))^2]}_{(2)}$$

$$\geq \underbrace{E[(Y - E[Y|X])^2]}_{\text{by } g(x) \stackrel{\text{(a.s.)}}{=} E[Y|X]}$$

# Proof

# Proof

# Orthogonality Property of MMSE

$$\begin{aligned} E(X) &= 0 \\ \Rightarrow X &\stackrel{\text{a.s.}}{=} 0 \end{aligned}$$

$$\text{Var}(X) = 0$$

$$E(X) > a \cdot p(X=a) > 0$$

$$E[X^2] = 0 \Rightarrow X = 0 \quad \text{a.s.}$$

$$0 \leq \text{Var}(X) = E[X^2] - E^2(X) \leq E[X^2]$$

$$\text{Var}(X) = 0$$

$$\Rightarrow E^2(X) \leq E[X^2] = 0$$

Theorem

$$\int_{B(a)} dx > 0$$

$E[Y|X]$  is a projector

(a) For any function  $\phi(\cdot)$ , one has

$$E[Y|X]$$

$$(Y - E[Y|X]) \perp \phi(x)$$

$$E[(Y - E[Y|X])\phi(X)] = 0 \quad \checkmark \quad (\#1)$$

(b) Moreover, if the function  $g(X)$  is such that

$$E[(Y - g(X))\phi(X)] = 0, \forall \phi(\cdot).$$

projector

$$\text{then } g(X) = \underline{E(Y|X)}$$

(#2)

is unique.

Proof  $E\left[\overbrace{(g(x) - E[Y|X])}^{\phi(x)}\right]^2 = E\left[\overbrace{(g(x) - E[Y|X])}^{\frac{g(x)-Y}{\square}}\right] \overbrace{(g(x) - E[Y|X])}^{\frac{Y-E[Y|X]}{\square}}$

$= E\left[\overbrace{(g(x) - E[Y|X])}^{\phi(x)}(g(x) - Y)\right] \textcircled{O} (\ast_2)$

$+ E\left[\overbrace{(g(x) - E[Y|X])}^{\phi(x)}\right] \overbrace{(Y - E[Y|X])}^{\textcircled{O} (\ast_1)} = \textcircled{O}$

$$\Rightarrow \underline{g(x) = E[Y|X]}$$

# Proof

# MMSE for Jointly Gaussian(Normal) Random Variables

$X, Y$  joint  $\Leftrightarrow \underline{a_1 X + a_2 Y \sim \text{Normal}}$

1°.  $Y - L[Y|X] \perp X$ .

2°.  $E[Y - L[Y|X]] = 0$

3°. Orthogonal  $\Leftrightarrow$  uncorrelated.

## Theorem

Normal

Let  $X, Y$  be jointly Gaussian random variables. Then

$$\underline{E[Y|X]} = \underline{L[Y|X]} = \underline{E(Y)} + \frac{\text{Cov}(X, Y)}{\text{Var}(X)}(\underline{X} - \underline{E(X)}).$$

4°.  $\underline{Y - L[Y|X]}$  and  $\underline{X}$  are joint Normal ✓

$$\begin{aligned} a_3(Y - \underline{L[Y|X]}) + a_4 X &= a_3(Y - a_5 X - a_6) + a_4 X \\ &= a_3 Y + (a_4 - a_5 a_3) X \end{aligned}$$

fasab

## Proof

5°.  $Y - L(Y|X)$  and  $X$  are independent.

6°.  $Y - L(Y|X)$  and  $\phi(x)$  are independent.

6°.  $\underline{Y - L(Y|X)}$  and  $\phi(x)$  are uncorrelated.

$$\underline{E[Y - L(Y|X)]} = 0$$

7°.  $\underline{Y - L(Y|X)} \perp \phi(x) \rightarrow \text{⑥φ(.)}$   
↓ projector.

8°.  $\underline{L(Y|X)} = \text{MSE}$

# Proof

## Example: Biased Coin Problem

1<sup>o</sup>. MMSE.  $E[\Theta|X=k]$  :  $\Theta \sim \text{Unif}(0,1)$

$X \sim \text{Beta}(a,b)$

$E[X] = \frac{a}{a+b}$

$= \text{Beta}(1,1)$

# of heads  $X|\Theta=\theta \sim \text{Bin}(n,\theta)$

By Beta-Binomial Conjugacy,  $\Theta|X=k \sim \text{Beta}(1+k, n-k)$

We wish to estimate the probability of landing heads, denoted by  $\theta$ , of a biased coin. We model  $\theta$  as the value of a random variable  $\Theta$  with a known prior PDF  $f_\Theta \sim \text{Unif}(0,1)$ . We consider  $n$  independent tosses and let  $X$  be the number of heads observed. Find the MMSE  $E(\Theta|X)$  and LLSE  $L(\Theta|X)$ .

$$E[\Theta|X=k] = \frac{k+1}{1+k+n-k} = \frac{k+1}{n+2}$$

$$\Rightarrow E[\Theta|X] = \frac{X+1}{n+2}$$

Solution 2<sup>o</sup>. LLSE.  $L(\theta | X) = E[\theta] + \frac{\text{Cov}(\theta, X)}{\text{Var}(X)}(X - E[X])$

$$\theta \sim \text{Unif}(0, 1) \Rightarrow E[\theta] = \frac{1}{2}, \quad \text{Var}(\theta) = \frac{1}{12} - E[\theta^2] = \frac{1}{3}$$

$$X|\theta = \theta \sim \text{Bin}(n, \theta) \Rightarrow E[X|\theta = \theta] = n\theta. \Rightarrow E[X|\theta] = n\theta.$$

$$\Rightarrow E[X] \stackrel{\text{Adam}}{=} E[E[X|\theta]] = E[n\theta] = nE[\theta] = \frac{1}{2}n$$

$$\text{Var}(X) \stackrel{\text{Eve}}{=} E[\text{Var}(X|\theta)] + \text{Var}[E[X|\theta]]$$

$$= E[n\theta(1-\theta)] + \text{Var}[n\theta]$$

$$= n[E(\theta) - E[\theta^2]] + n^2 \text{Var}(\theta)$$

$$= n\left[\frac{1}{2} - \frac{1}{3}\right] + n^2 \cdot \frac{1}{12}$$

$$= \frac{n}{12}(n+2)$$

Solution

$$\begin{aligned}
 \text{Cov}(x, \theta) &= \overline{E[\theta x]} - \overline{E[\theta]} \cdot \overline{E[x]} \\
 &= \overline{E[E[\theta x | \theta]]} - \frac{1}{2} \cdot \frac{n}{2} \\
 &= E[\overline{\theta E[x | \theta]}] - \frac{n}{4} \\
 &= E[\theta \cdot n\theta] - \frac{n}{4} \\
 &= n E[\theta^2] - \frac{n}{4} \\
 &= n \cdot \frac{1}{3} - \frac{n}{4} = \frac{1}{12}n
 \end{aligned}$$

$\Rightarrow$  LLSE

$$\begin{aligned}
 L(\theta | x) &= E[\theta] + \frac{\text{Cov}(\theta, x)}{\text{Var}(x)}(x - E(x)) \\
 &= \frac{1}{2} + \frac{\frac{1}{12}n}{\frac{n}{12}(n+2)}(x - \frac{n}{2}) \\
 &= \frac{x+1}{n+2} = E[\theta | x]
 \end{aligned}$$

# Solution

# Remark: Statistical Learning Perspective

$$E[Y|X]$$

$$E[Y|X=x]$$

Regression function

- In general, MMSE is a highly nonlinear function.
- Adoption of various approximation methods leads to various learning methods
  - ▶ Linear regression
  - ▶ Logistic regression
  - ▶ Polynomial regression
  - ▶ Regression with Spline functions
  - ▶ Neural network

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# 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...

# History



# 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

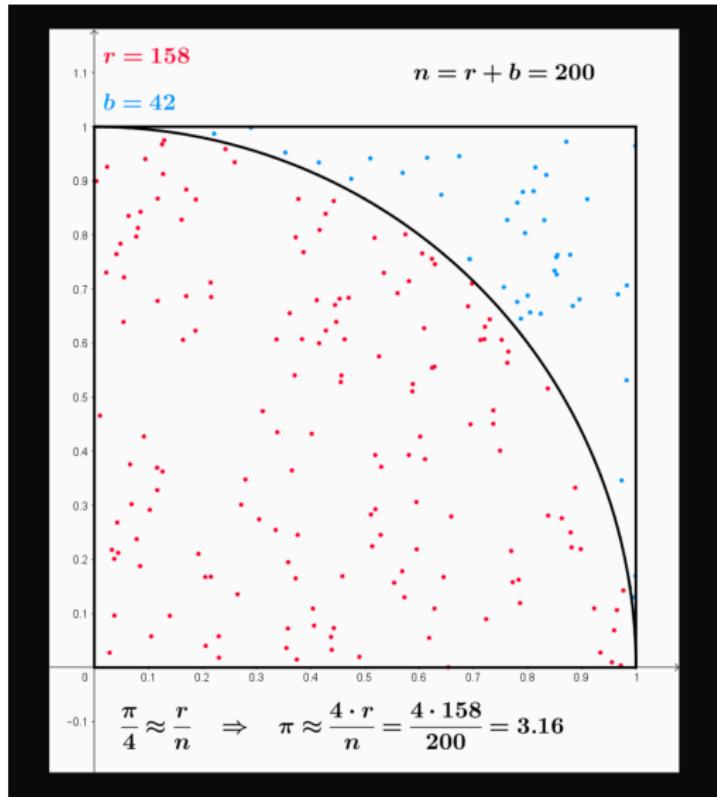
$$c = \overbrace{E_{\pi}(h(x))}^{\text{red circle}} = \overbrace{\int f(x) h(x) dx}^{\text{red circle}}$$

- Optimizing a target function to find its maxima or minima

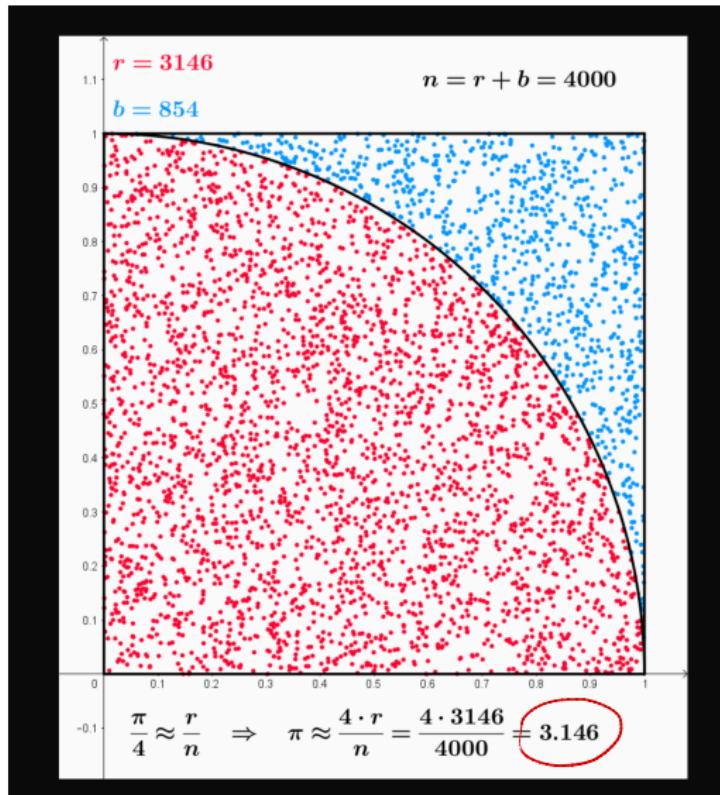
# Two Main Categories

- Sequential Monte Carlo(Our Focus in this Lecture): preserves & propagates a population of examples by sequential sampling and importance reweighing, often in a low dimensional state space.
- Markov Chain Monte Carlo(Next Lecture): simulate a Markov chain to explore the state space with a stationary probability designed to converge to a given target probability

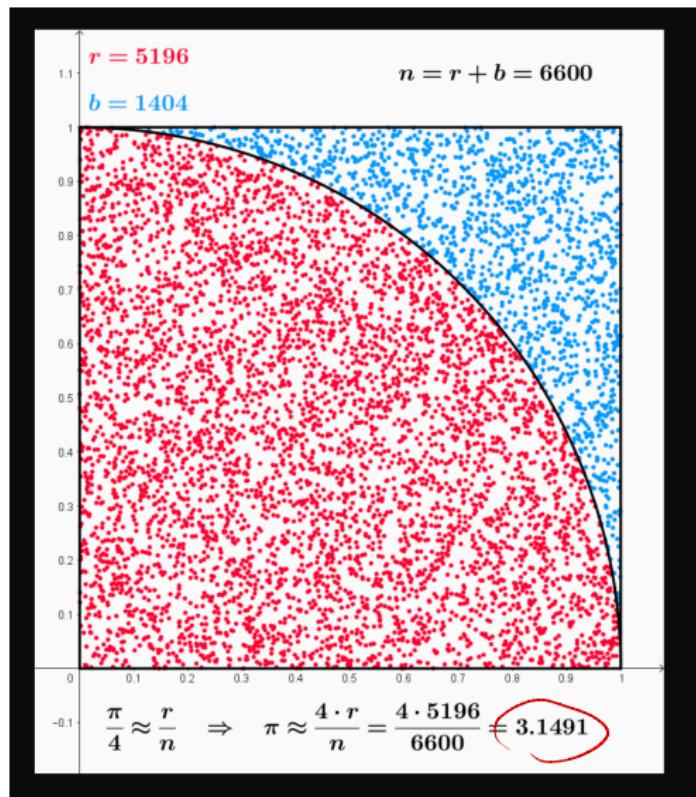
# Classical Example: Estimation of $\pi$



# Classical Example: Estimation of $\pi$



# Classical Example: Estimation of $\pi$

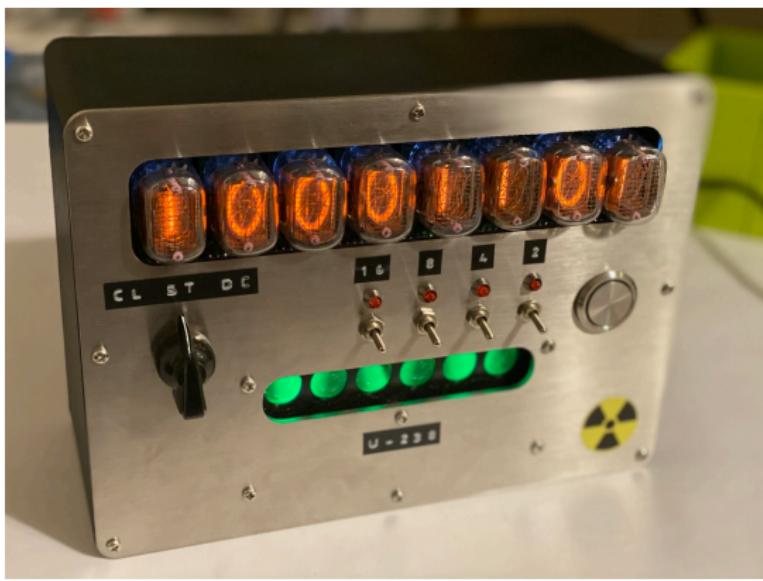


# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# 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)
- Current days: computer simulation with deterministic algorithms, also called pseudorandom number generator

# Sampling

- Assuming an algorithm is available for generating 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

$$P(F^{-1}(U) \leq x) = P(U \leq F(x)) = F(x)$$

## 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  $F$ .

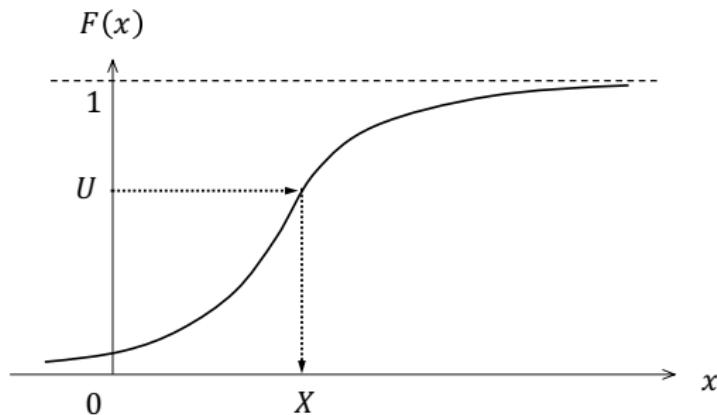
**output:** Random variable  $X$  distributed according to  $F$ .

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

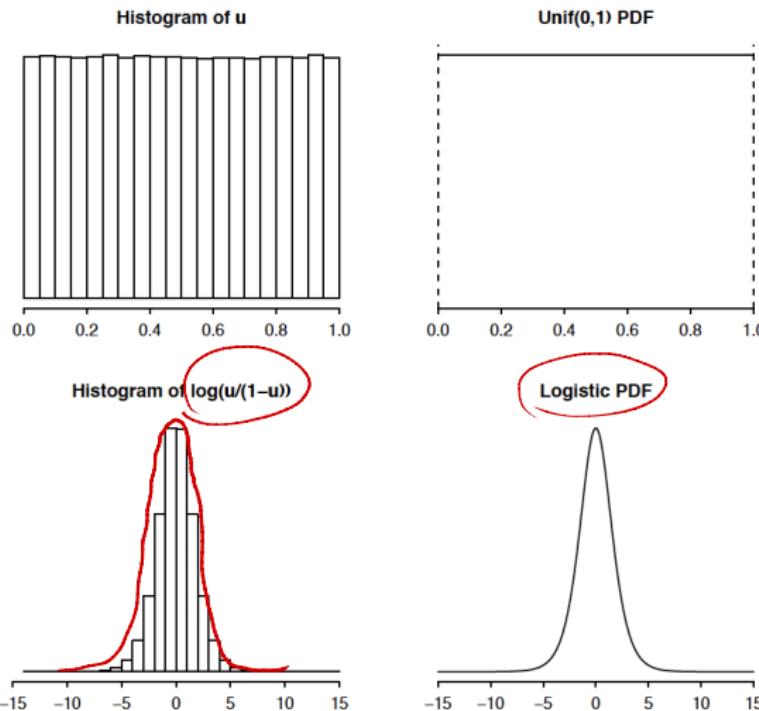
2:  $X \leftarrow F^{-1}(U)$

3: **return**  $X$

---



# Histogram & PDF: Example



## Box-Muller Method: Recall

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.

---

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

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

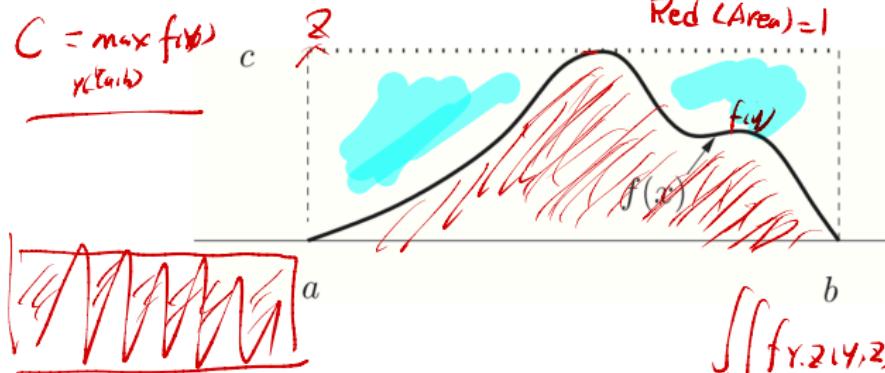
- 1: Generate two independent random variables,  $U_1$  and  $U_2$ , from  $\text{Unif}(0, 1)$ .
- 2:  $X \leftarrow (-2 \ln U_1)^{1/2} \cos(2\pi U_2)$        $T \sim \text{Expo}(1)$
- 3:  $Y \leftarrow (-2 \ln U_1)^{1/2} \sin(2\pi U_2)$        $F_T(t) = 1 - e^{-t}, t \geq 0$
- 4: **return**  $X, Y$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

$$U_1 \sim \text{Unif}(0,1)$$

$$U_2 \sim \text{Unif}(0,1)$$

# Acceptance-Rejection Method

$$C = \max_{y \in [a, b]} f(y)$$



① PDF  $f(y) = \frac{f(y)}{\int_0^b f(y) dy}$

$$= \int_a^b f(y) f_{Y,Z}(y,z) dy$$

$$\frac{f_{Y,Z}(y,z)}{f_{Y,Z}(y,z)} = 1$$

$$= \frac{1}{\text{Red(Area)}}$$

$$\int \int f_{Y,Z}(y,z) dy dz$$

## Algorithm Acceptance-Rejection Algorithm

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

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

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

$$= \frac{1}{\text{Red(Area)}} \left( \int_a^b f(y) dy \right)$$

# Acceptance-Rejection Method

② Uniform Sampling over Triangle

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

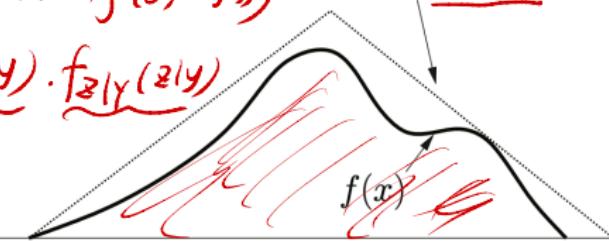
$$\Rightarrow f_{Y,Z}(y,z) = f_Y(y) \cdot f_Z(z|y)$$

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

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

$$\textcircled{1} \quad x \in [a, b], \quad \phi(x) \geq f(x) \Rightarrow Cg(x) \geq f(x) \\ \text{Valid g: PDF} \\ \Rightarrow C \geq \left[ \sup_x \frac{f(x)}{g(x)} \right]$$

$$\phi(x) = Cg(x)$$



(Area(Triangle))

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

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

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

$$= C$$

## Algorithm Acceptance-Rejection Algorithm

Step 1: Generate  $Y \sim g$ .

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

Step 3: If  $Z \leq f(Y)$ , 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

- 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:

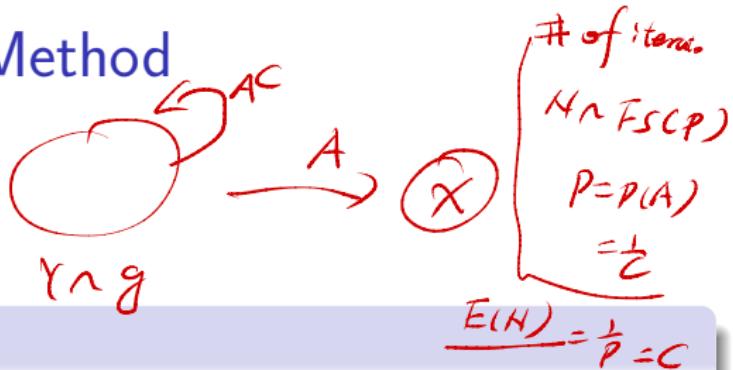
Step 1: Generate  $\underline{Y \sim g}$ .

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

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

---

# Acceptance-Rejection Method



## 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  $A = \{U \leq \frac{f(Y)}{c \cdot g(Y)}\}$ .  $U$  uniform

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

$$\begin{aligned} 1^{\circ} \quad P(A|Y=y) &= P(U \leq \frac{f(y)}{c \cdot g(y)} | Y=y) = P(U \leq \frac{f(y)}{c \cdot g(y)} | Y=y) \\ &= P(U \leq \frac{f(y)}{c \cdot g(y)}) \stackrel{(0,1)}{=} \frac{f(y)}{c \cdot g(y)} \end{aligned}$$

$c > \sup_y \frac{f(y)}{g(y)}$

$$\begin{aligned} 2^{\circ} \quad P(A) &\stackrel{\text{LorP}}{=} \int P(A|Y=y) \cdot f_Y(y) dy \quad (\text{Y~\&~} \mathbb{R}) \\ &= \int \frac{f(y)}{c \cdot g(y)} \cdot \cancel{dy} dy = \frac{1}{c} \int \cancel{f(y)} dy = \underline{\frac{1}{c}} \leq 1 \end{aligned}$$

(C.1)

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

$$\underline{f_X(y)} = f(y)$$

# Proof

## Example: Beta Distribution

- 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  $X \sim \text{Beta}(a, b)$ .

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

Solution

$$\text{object PDF } f(x) = 20x(1-x)^3, 0 < x < 1$$

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

$$C \geq \sup_{y \in (0,1)} \frac{f(y)}{g(y)} = \sup_{y \in (0,1)} \frac{20y(1-y)^3}{1} = \sup_{y \in (0,1)} 20y(1-y)^3 \Rightarrow y^* = \frac{1}{4}$$

$$\Rightarrow C \geq \frac{f(y^*)}{g(y^*)} = \frac{135}{64} > 1, \Rightarrow \text{choose } C = \frac{135}{64}$$

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

---

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

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

3: If  $U \leq \frac{f(Y)}{C \cdot g(Y)} = \frac{256}{135} Y(1-Y)^3$ ,  $X=Y$ .

Otherwise Reject  $Y$ , Go back to Step 1.

# Solution

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# 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**

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

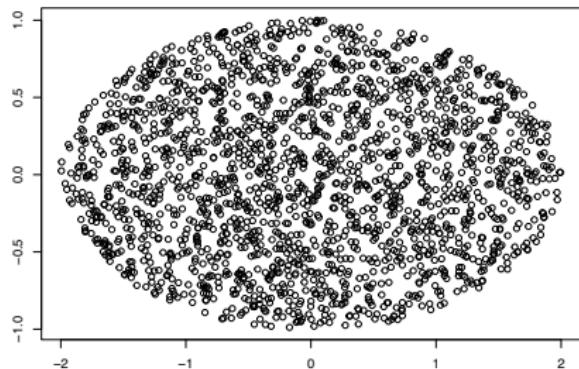
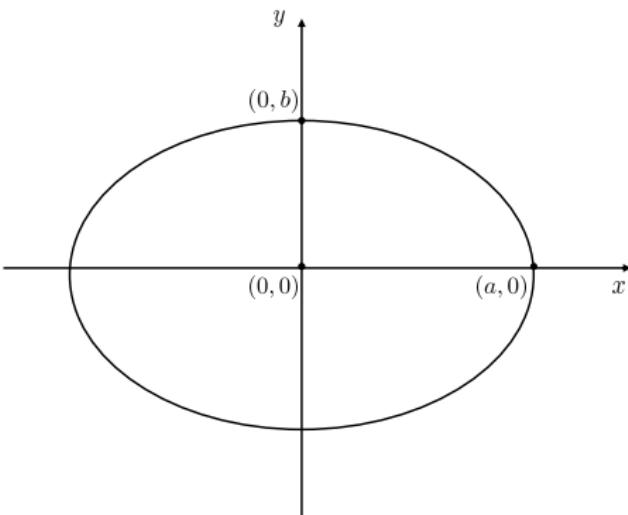
Also assume that the determinant of the Jacobian matrix is never 0. Then the joint PDF of  $\mathbf{Y}$  is

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

# Example: Generate Uniform Distribution over An Ellipse

- Ellipse:

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



Solution ① Object PDF  $f_{X,Y}(x,y) = \frac{1}{\pi ab}$ ,  $\forall x,y \in E_{(a,b)}$

$$X = p \cos \theta$$

$$Y = p b \cdot \sin \theta$$

$(p \in [0,1], \theta \in [0,2\pi])$

② Jacob Matrix  $\begin{bmatrix} \frac{\partial X}{\partial p} & \frac{\partial X}{\partial \theta} \\ \frac{\partial Y}{\partial p} & \frac{\partial Y}{\partial \theta} \end{bmatrix} = \begin{bmatrix} \cos \theta & -p \sin \theta \\ b \sin \theta & pb \cos \theta \end{bmatrix}$

$J = \det(\ ) = pab.$

$$\Rightarrow f_{R,\Theta}(p,\theta) = f_{X,Y}(x,y) \cdot |J| = \frac{1}{\pi ab} \cdot pab = \frac{p}{\pi}.$$

$$\Rightarrow f_R(p) = \int_0^{2\pi} \frac{p}{\pi} d\theta = 2p, 0 \leq p \leq 1 \Rightarrow \underline{F_R(p) = p^2, 0 \leq p \leq 1}. \quad \Rightarrow F_R^{-1}(z) = \sqrt{z}, 0 \leq z \leq 1.$$

$$f_\Theta(\theta) = \int_0^1 \frac{p}{\pi} dp = \frac{1}{2\pi}, 0 \leq \theta \leq 2\pi, \Rightarrow \underline{\Theta \sim \text{unif}(0,2\pi)} = 2\pi \text{unif}(0,1)$$

$f_{R,\Theta}(p,\theta) = f_R(p) f_\Theta(\theta)$ , R and  $\Theta$  are independent.

## Solution

(3)

$$X = \rho a \cos \theta, \quad Y = \rho b \sin \theta$$

$U_1, U_2$  independent uniform

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

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

$(X, Y)$

# Solution

# Change of Variables

$$\begin{aligned} A &\text{ is } n \times n \text{ matrix} \\ \det(A^T A) &= \det^2(A) \end{aligned}$$

## 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(Surface)

- Sphere Surface:

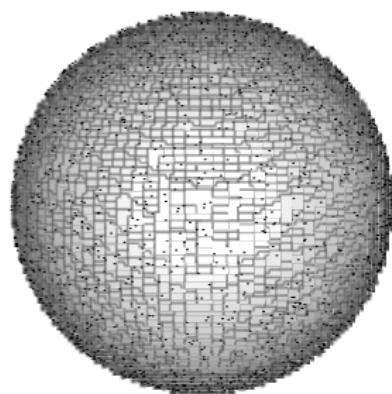
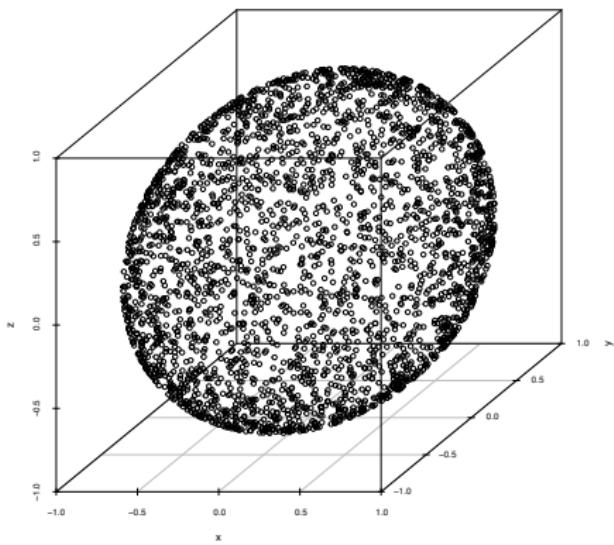
$$S_2(r) = \{(x, y, z) \in \mathbb{R}^3 : x^2 + y^2 + z^2 = r^2\}.$$

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

(Ball)

$B^2 r^3$

$4\pi r^2$



Solution ①  $f_{x,y,z}(x,y,z) = \frac{1}{4\pi r^2}, \quad (x,y,z) \in S_2(r)$

$$x = rs \cos \theta \cos \phi, \quad \theta \in [0, \pi]$$

$$y = rs \sin \theta \sin \phi, \quad \phi \in [0, 2\pi]$$

$$z = r \cos \theta$$

$$(x, y, z) \rightarrow (r, \theta, \phi)$$

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

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

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

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

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

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

$$\Rightarrow f_{\theta, \phi}(\theta, \phi) = f_{\theta}(\theta) f_{\phi}(\phi). \Rightarrow \theta \text{ and } \phi \text{ are independent.}$$

$\phi \sim \text{unif}(0, 2\pi) \Rightarrow \text{unif}(0, 1)$

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

---

$$\text{if } \theta = \arccos(1 - 2s) \Rightarrow \cos \theta = \underline{1 - 2s} \Rightarrow \sin^2 \theta = 1 - \cos^2 \theta \\ = 4s(1-s).$$

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

## Solution

$$\textcircled{3} \quad X = r \sin \theta \cos \phi = r \cdot 2\sqrt{1-s^2} \cos(2\pi s)$$

$$Y = r \sin \theta \sin \phi = r \cdot 2\sqrt{1-s^2} \sin(2\pi s)$$

$$Z = r \cos \theta = r(1-2s)$$

---

\textcircled{4} .  $U_1, U_2$  independently unif(0,1)

$$X \leftarrow r \cdot 2\sqrt{U_1(1-U_1)} \cdot \cos(2\pi U_2)$$

$$Y \leftarrow r \cdot 2\sqrt{U_1(1-U_1)} \cdot \sin(2\pi U_2)$$

$$Z \leftarrow r(1-2U_1)$$

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# Monte Carlo Integration

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

$$\underline{E[g(X)]} \approx \frac{1}{n} \sum_{i=1}^n g(X_i).$$

$E[g(X)]$   
 $X \sim \text{Unif}(a, b)$

- Now we have integration

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

$f(x)$  POF  
 $X \sim \text{Unif}(a, b)$

- Drawing  $n$  samples (empirical samples) from  $\text{Unif}(a, b)$ :

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

- Monte Carlo Integration:

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

# Example: $\pi$ as An Integration

Evaluate the integration

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

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

$$\int_0^1 \frac{4}{1+x^2} dx \approx \frac{1}{n} \sum_{i=1}^n \frac{4}{1+X_i^2}.$$

## Example

Evaluate the integration

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

- Corresponding

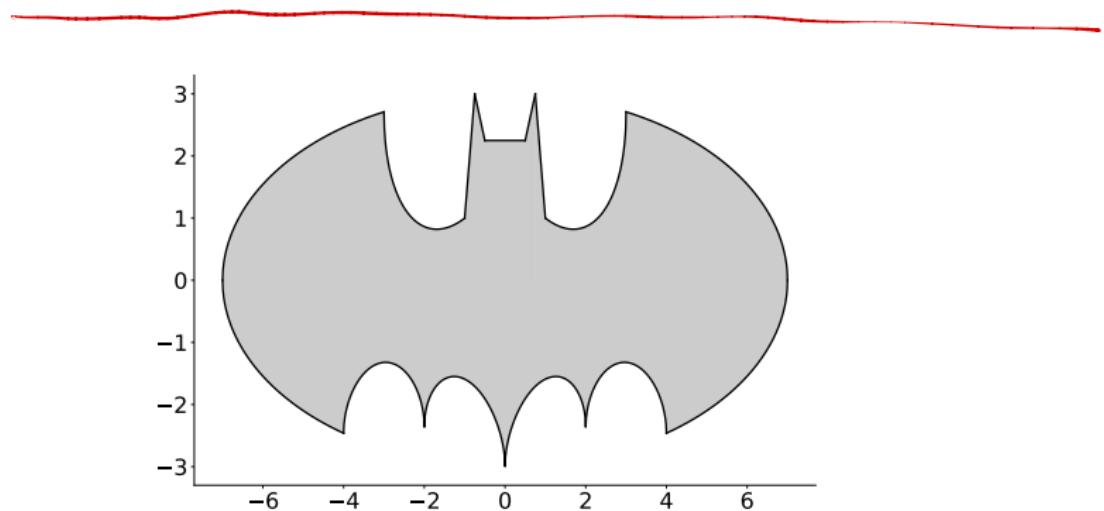
$$g(x) = \sqrt{x + \sqrt{x + \sqrt{x + \sqrt{x}}}}$$

- $X_1, \dots, X_n$ : samples from Unif(0, 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>



# 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.

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

# Importance Sampling

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

- Hence, given an iid sample  $Y_1, \dots, Y_n$  from PDF  $g$ , 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

$$C = P(Y > 8) = E[\underline{I(Y > 8)}] \approx \frac{1}{n} \sum_{j=1}^n I(Y_j > 8)$$

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

choose  $g \sim N(8, 1)$ ,  $Y_1, \dots, Y_n \sim g$   $f \sim N(0, 1)$

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

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

$$n = 50000 ; C \approx 6.25 \times 10^{-16}$$

# Solution

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

# 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  $\sigma^2/n$ .

# Asymptotic Analysis: Strong Law of Large Numbers (SLLN)

$$n \rightarrow \infty$$

## 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.*

# Asymptotic Analysis: Weak Law of Large Numbers (WLLN)

## 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).

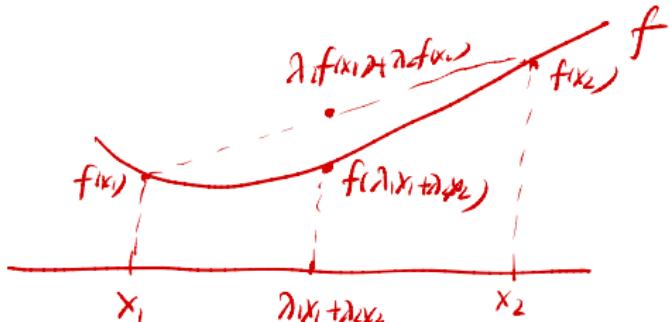
## Finite Sample Analysis: Inequalities

Theorem (Cauchy-Schwarz Inequality)

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

2<sup>o</sup>. Subgradient  $\underline{g^0(\cdot)}$

$$\frac{g(y) \geq g(x) + \underline{(g^0(x))(y-x)}}{y \in X}$$

Theorem  $\underline{X \subseteq E[X]}$

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.

$$\begin{aligned} 1^o. \quad & g(y) \geq g(x) + g'(x)(y-x) \\ & \underset{y \notin X}{\text{---}} \quad \underset{x \in E[X]}{\text{---}} \Rightarrow g(x) \geq \underline{g(E[X])} + g'(E[X])(x - \underline{E[X]}) \\ & \Rightarrow \underline{E[g(x)]} \geq g(E[X]) + g'(E[X])\underline{E[x - E[X]]} \\ & \quad \quad \quad = g(E[X]) \end{aligned}$$

## Quick Examples

$g$  is convex  
Concave

$$E[g(x)] \geq g[E(x)],$$
$$\leq$$

$$g''(\cdot) \geq 0$$
$$g''(\cdot) \leq 0$$

1<sup>o</sup>.  $g(x) = x^2, x \in R$ , convex  $\Rightarrow E[x^2] \geq (E[x])^2$  ✓  
 $\text{Var}(x) = E(x^2) - E^2(x) \geq 0$

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  $\Rightarrow E[\log x] \leq \underline{\log(E[x])}$

# 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 r.v.  $Y$  s.t.

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

$$\textcircled{2} \quad H(X) \triangleq \sum_{j=1}^n p_j \log_2 \frac{1}{p_j} = \underline{E[\log_2 Y]} \leq \log_2 E(Y) = \log_2 n$$

$$\forall p_1, p_2, \dots, p_n \quad \underbrace{p_1 + \dots + p_n = 1}_{\text{P.F.}} \Rightarrow \max_{p_1, \dots, p_n} H(X) \leq \log_2 n$$

$$\textcircled{3} \quad \text{when } X \sim \text{unif}(\frac{1}{n}), \quad \underbrace{p_1 = p_2 = \dots = p_n = \frac{1}{n}}, \quad H(X) = \sum_{j=1}^n \frac{1}{n} \log_2 \frac{1}{n}$$

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

## Kullback-Leibler Divergence

PDF:  $p(x)$ ,  $r(x)$ .  $x \in X$ .

$$\begin{aligned} D(p, r) &= \int p(x) \log_2 \frac{p(x)}{r(x)} dx \\ &= - \int p(x) \log_2 \frac{r(x)}{p(x)} dx \end{aligned}$$

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

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

Show that the Kullback-Leibler divergence is nonnegative.

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

Construct a r.v.  $Y$ . s.t.

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

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

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

(2)  $\underbrace{\log_2 X \leq x-1}$

$$D(p, r) \geq -\sum_{j=1}^n p_j \left( \frac{r_j}{p_j} - 1 \right)$$

$$= -\sum_{j=1}^n (r_j - p_j) = -\sum_{j=1}^n r_j + \sum_{j=1}^n p_j \\ = -1 + 1 = 0$$

# Markov's Inequality

$$P(|X - E(X)| \geq a)$$

$a \uparrow$

prob  $\downarrow$

Concentration

Inequality

$$\left\{ \begin{array}{l} \frac{1}{a} \\ \frac{1}{a^2} \\ e^{-a} \\ e^{-a} \end{array} \right.$$

## Theorem

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

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

Proof  $P(|x| \geq a) \leq \frac{1}{a} E(|x|), a > 0.$

①  $Y = \frac{1}{a} |x| \geq 0 \quad , \quad \underbrace{I(Y \geq 1)}_{0 \leq Y < 1} \leq Y.$  LHS  $|Y \geq 1|$   $1 \leq Y$  RHS  $0 \leq Y$

②  $E[I(Y \geq 1)] \leq E[Y]$

$$\underbrace{P(Y \geq 1)}_{\text{by}} \leq \underbrace{E[Y]}_{= E[\frac{1}{a}|x|] = \frac{1}{a} E(|x|)}$$

$$P(\frac{1}{a}|x| \geq 1) \leq$$

↙

$$\underbrace{P(|x| \geq a)}_{\text{↙}}$$

# Chebyshev's Inequality

Markov's inequality

$$\begin{aligned} P(|X-\mu| \geq a) &= P(|X-\mu|^2 \geq a^2) \leq \frac{1}{a^2} E(|X-\mu|^2) \\ &= \frac{1}{a^2} \text{Var}(X) \\ &= \frac{1}{a^2} \sigma^2 \end{aligned}$$

## Theorem

Let  $X$  have mean  $\mu$  and variance  $\sigma^2$ . Then for any  $a > 0$ ,

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

# Proof

## Chernoff's Inequality

$t > 0,$

$$P(X \geq a) = P(tx \geq ta)$$

$$= P(e^{tX} \geq e^{ta})$$

Moment's inequality  
 $\leq$

$$\frac{E(e^{tx})}{e^{ta}} = f(t)$$

### Theorem

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

$$P(X \geq a) \leq \frac{E(e^{tX})}{e^{ta}}.$$

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

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

# Proof

# Chernoff's Technique

$$\Pr[X \geq a]$$

$$= \Pr[e^{tX} \geq e^{ta}]$$

$$= \Pr[e^{tX} \geq e^{ta}]$$

$$\leq \frac{E[e^{tX}]}{e^{ta}}.$$

## 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

# Hoeffding Lemma

$$X \sim N(\mu, \sigma^2)$$
$$E[e^{tX}] = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$$
$$\mu t + \frac{1}{2}\sigma^2 t^2 = e^{\frac{1}{2}\sigma^2 t^2}$$

## Lemma

Let the random variable  $X$  satisfy  $\mathbb{E}(X) = 0$  and  $a \leq X \leq b$ , where  $a$  and  $b$  are constants. Then for any  $\lambda > 0$ ,

$$\mathbb{E}(e^{\lambda X}) \leq e^{\frac{1}{8}\lambda^2(b-a)^2}.$$

Sub-Gaussian

# Useful Analysis Tools

- Jensen's inequality: if  $f$  is convex,  $0 \leq \lambda_1, \lambda_2 \leq 1, \lambda_1 + \lambda_2 = 1$ , then for any  $x_1, x_2$ ,

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

- Taylor's theorem or Taylor's expansion: If all the derivatives of a function  $f(x)$  exist at point  $a$ , then for any positive integer  $k$ , there exist a real number  $\theta$  between  $a$  and  $x$  such that

$$f(x) = f(a) + \cdots + \frac{f^{(k)}(a)}{k!}(x-a)^k + \frac{f^{(k+1)}(\theta)}{(k+1)!}(x-a)^{k+1}.$$

Hoeffding Bound

$$(1) \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i - \mu = \frac{1}{n} \sum_{i=1}^n Z_i$$

$$Z_i = X_i - \mu = X_i - E[X_i] \Rightarrow E[Z_i] = 0, i=1,2,\dots,n.$$

$$Z_1, \dots, Z_n \text{ independent} \quad E[\bar{X}] = 0.$$

## 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  $\epsilon \geq 0$ ,

$$\{|\bar{X}| \geq \epsilon\} \subset \{\bar{X} > \epsilon\} + \{\bar{X} < -\epsilon\}$$

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

$O(e^{-n\epsilon^2})$

$$(2) \quad \forall \lambda > 0, \quad \underline{P(\bar{X} \geq \epsilon)} = P(e^{\lambda \bar{X}} \geq e^{\lambda \epsilon})$$

$$\leq \frac{1}{e^{\lambda \epsilon}} \underline{E[e^{\lambda \bar{X}}]} = \frac{1}{e^{\lambda \epsilon}} \underline{E\left[e^{\lambda \cdot \frac{1}{n} \sum_{i=1}^n Z_i}\right]}$$

Proof

$$\forall \lambda > 0 ; P(Z \geq \varepsilon) \leq e^{-\lambda \varepsilon} \underbrace{\prod_{i=1}^n E[e^{\lambda \frac{1}{n} Z_i}]}_{E[e^{\lambda \frac{1}{n} Z_i}]}$$

$$\left( \frac{1}{n} Z_i \right) := E\left(\frac{1}{n} Z_i\right) = \frac{1}{n} E(Z_i) = 0$$

$$a \leq X_i \leq b \Rightarrow a - \mu \leq \underline{X_i - \mu} \leq b - \mu$$

$$\Rightarrow \underline{\frac{1}{n}(a-\mu)} \leq \frac{1}{n} Z_i \leq \underline{\frac{1}{n}(b-\mu)}$$

By Hoeffding's Lemma .  $E[e^{\lambda \cdot \frac{1}{n} Z_i}] \leq e^{\frac{1}{8} \lambda^2 \left( \frac{1}{n}(b-\mu) - \frac{1}{n}(a-\mu) \right)^2}$

$$(3) \quad \underline{P(Z \geq \varepsilon)} \leq e^{-\lambda \varepsilon} \cdot \overbrace{\prod_{i=1}^n \left[ e^{\frac{1}{8} \lambda^2 \frac{(b-a)^2}{n^2}} \right]}^{= e^{-\lambda \varepsilon} \cdot e^{\frac{1}{8} \lambda^2 \frac{(b-a)^2}{n}}} , \forall i = 1, 2, \dots, n$$
$$= e^{-\lambda \varepsilon + \frac{1}{8} \lambda^2 \frac{(b-a)^2}{n}} = e^{\underline{g(\lambda)}} \quad \text{if } \lambda > 0$$

$$g(\lambda) = -\lambda \varepsilon + \frac{1}{8} \lambda^2 \frac{(b-a)^2}{n} \Rightarrow g'(\lambda^*) = 0 ; g''(\lambda^*) > 0.$$

$$\Rightarrow \lambda^* = \frac{4n\varepsilon}{(b-a)^2} \Rightarrow g(\lambda^*) = \frac{-2n\varepsilon^2}{(b-a)^2}$$



Proof  $\underline{P(Z \geq \varepsilon)} \leq \inf_{\lambda > 0} e^{g(\lambda)} = e^{g(\lambda^*)} = e^{-\frac{2n\varepsilon^2}{(b-a)^2}}$

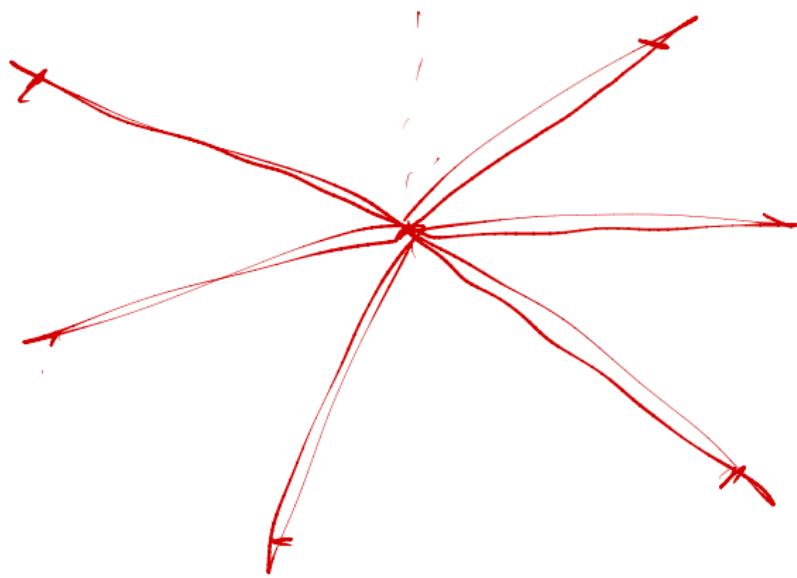
(4)  $\underline{P(Z \leq -\varepsilon)} \leq e^{-\frac{2n\varepsilon^2}{(b-a)^2}}$

$$\Rightarrow P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - \mu\right| \geq \varepsilon\right) = P(|Z| \geq \varepsilon)$$

$$= P(Z \geq \varepsilon) + P(Z \leq -\varepsilon) = 2e^{-\frac{2n\varepsilon^2}{(b-a)^2}}$$

$n \rightarrow \infty \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \mu$

# Proof



# More General Hoeffding Bound

## Theorem

Let the random variables  $X_1, X_2, \dots, X_n$  be independent, with  $a_k \leq X_k \leq b_k$  for each  $k$ , where  $a_k, b_k$  are constants. Let  $S_n = \sum_{k=1}^n X_k$  and let  $\mu = \mathbb{E}(S_n)$ . Then for any  $t \geq 0$ ,

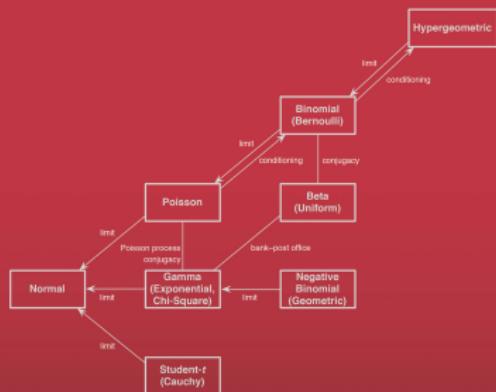
$$\mathbb{P}(|S_n - \mu| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{k=1}^n (b_k - a_k)^2}}.$$

# Outline

- 1 Stochastic Laws in Random Worlds
- 2 Statistical Inference
- 3 Conjugate Prior: A Weapon of Bayesian
- 4 Conditional Expectation
- 5 Monte Carlo Method: Another Weapon of Bayesian
- 6 Sampling: Random Variable Generation
- 7 Sampling: Random Vector Generation
- 8 Monte Carlo Integration
- 9 Performance Analysis of Monte Carlo Integration
- 10 References

Texts in Statistical Science

# Introduction to Probability



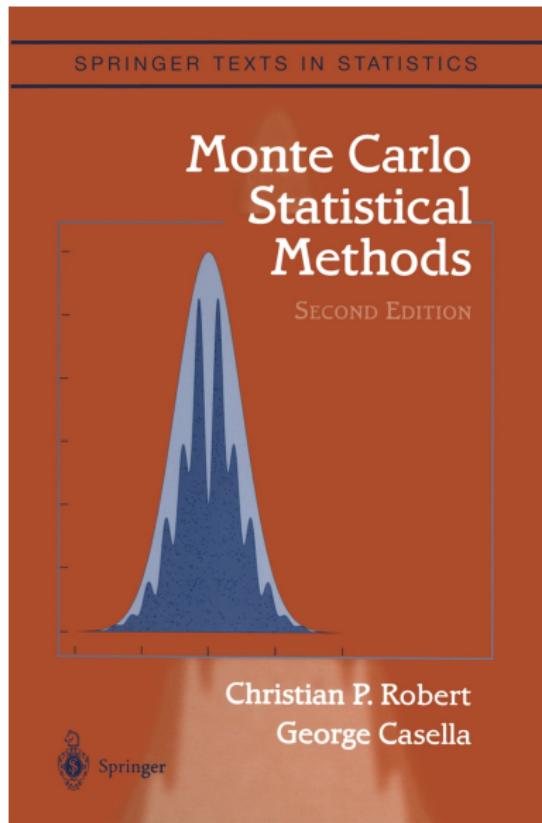
Joseph K. Blitzstein  
Jessica Hwang



CRC Press  
Taylor & Francis Group  
A CHAPMAN & HALL BOOK

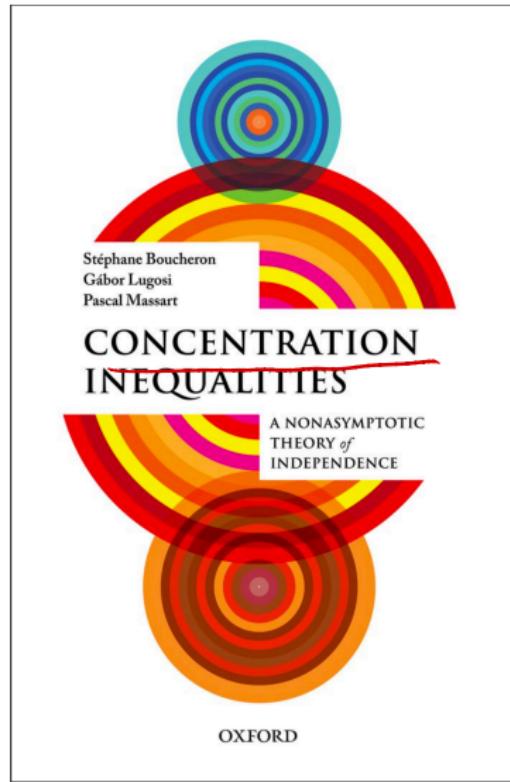
Joseph K. Blitzstein &  
Jessica Hwang

- Introduction to Probability
- Chapman & Hall/CRC, 2014.
- Chapman & Hall/CRC, 2019.
- Chapters 8 & 9



Christian P. Robert &  
George Casella

- Monte Carlo Statistical Methods
- Springer, 2010



## Stephane Boucheron & Gabor Lugosi & Pascal Massart

- Concentration Inequalities
- Oxford, 2013