

Hoff Ch.4 : Monte Carlo

- Previously. Conjugacy led to posterior distributions with simple, closed-form expressions for, e.g., mean and variance.

- Beta + Binomial = Beta
- Gamma + Poisson = Gamma

However, even for these simple conjugate models, posterior inference can be hard.

Looking back at the dental health CDC data:

- 48 40 year old men with 23 unhealthy teeth total
- 57 40 year old women w/ 33 unhealthy teeth total.
- Gamma(1,1) prior

$$\Rightarrow p(\theta_m | \vec{y}_m) = \text{Gamma}(24, 49)$$

$$p(\theta_w | \vec{y}_w) = \text{Gamma}(34, 58)$$

Probability that  $\Theta_w > \Theta_m$ :

$$\mathbb{P}(\Theta_w > \Theta_m | \vec{y}_m, \vec{y}_w) = \int_0^\infty \int_0^{\Theta_w} p(\Theta_m, \Theta_w | \vec{y}_m, \vec{y}_w) d\Theta_m d\Theta_w$$

$\nearrow$

$$= \int_0^\infty \int_0^{\Theta_w} p(\Theta_m | \vec{y}_m) p(\Theta_w | \vec{y}_w) d\Theta_m d\Theta_w$$

Assuming  $\Theta_m, \Theta_w$  conditionally independent given data

$$= \frac{49^{24} 58^{34}}{\Gamma(24)\Gamma(34)} \int_0^\infty \int_0^{\Theta_w} \Theta_m^{23} e^{-49\Theta_m} \Theta_w^{33} e^{-58\Theta_w} d\Theta_m d\Theta_w.$$

Problem. Closed-form expressions often unavailable for posterior quantities of interest.

Integrals of the form

$$\mathbb{E}[f(\theta) | y] = \int_{\Theta} f(\theta) p(\theta | y) d\theta$$

often extremely hard to compute exactly.

- Law of large numbers. If  $X_1, \dots, X_S$  are iid copies of r.v.  $X$ , then

$$\lim_{S \rightarrow \infty} \frac{1}{S} \sum_{i=1}^S X_i = E[X].$$

- More generally, for (almost) any function  $f$ ,

$$\lim_{S \rightarrow \infty} \frac{1}{S} \sum_{i=1}^S f(X_i) = E[f(X)].$$

Side note: For any  $\mathbb{Z}$ -valued random variable  $Z$  and any event  $A \subseteq \mathbb{Z}$ , we have

$$P(Z \in A) = E[1\!(Z \in A)]$$

where  $1\!(Z \in A)$  is the indicator function

$$1\!(Z \in A) = \begin{cases} 1 & \text{if } Z \in A \\ 0 & \text{if } Z \notin A \end{cases} .$$

Computing  $P(\Theta_w > \Theta_m | \vec{y}_m, \vec{y}_w)$ .

- Draw  $S$  samples  $\Theta_{m,i} \stackrel{\text{iid}}{\sim} \text{Gamma}(24, 49)$
- Draw  $S$  samples  $\Theta_{w,i} \stackrel{\text{iid}}{\sim} \text{Gamma}(34, 58)$
- $P(\Theta_w > \Theta_m | \vec{y}_m, \vec{y}_w) = E\left[1(\Theta_w > \Theta_m) | \vec{y}_m, \vec{y}_w\right]$   
 $\approx \frac{1}{S} \sum_{i=1}^S 1(\Theta_{w,i} > \Theta_{m,i}).$

## Monte Carlo: LLN for posterior inference

- Suppose we draw  $S$  iid samples

$$\theta_i \stackrel{\text{iid}}{\sim} p(\theta|y).$$

Then for large  $S$ , by LLN,

$$\frac{1}{S} \sum_{i=1}^S f(\theta_i) \approx \int_{\Theta} f(\theta) p(\theta|y) d\theta.$$

- Posterior mean:

$$\bar{\theta}_S = \frac{1}{S} \sum_{i=1}^S \theta_i \rightarrow \mathbb{E}[\theta | y]$$

- Posterior variance:

$$\frac{1}{S-1} \sum_{i=1}^S (\theta_i - \bar{\theta}_S)^2 \rightarrow \text{Var}[\theta | y]$$

- Probability  $\Theta \leq c$ :

$$\frac{1}{S} \sum_{i=1}^S \mathbb{1}(\Theta_i \leq c) \rightarrow \mathbb{E}[\mathbb{1}(\Theta \leq c) | y]$$
$$= P(\Theta \leq c | y).$$

- Quantiles.

$\alpha \in (0, 1)$

The  $\alpha$  quantile of  $\Theta | y$  is a value  $q_\alpha$  such that

$$P(\theta \leq q_\alpha | y) = \alpha .$$

This definition requires some assumptions about the distribution of  $X$ , but we'll skip the more general definition for simplicity.

To estimate  $q_\alpha$  using Monte Carlo:

(i) Draw  $\theta_i \stackrel{iid}{\sim} p(\theta | y)$ ,  $i = 1, \dots, S$

(ii) Pick  $q_\alpha(S)$  so that 100 $\alpha$  % of samples are  $\leq q_\alpha(S)$ .

$\Rightarrow q_\alpha(S) \longrightarrow q_\alpha \text{ as } S \longrightarrow \infty .$

## Bayesian stats in a nutshell.

- Come up with model
  - Prior, likelihood
- Compute posterior,  $p(\theta | y)$
- Draw samples  $\theta_1, \dots, \theta_S$  from  $p(\theta | y)$
- Make inferences via Monte Carlo.

## Monte Carlo error

- How much does  $\bar{X}_S = \frac{1}{S} \sum_{i=1}^S X_i$  fluctuate around its mean,  $E[\bar{X}_S] = E[X] := \mu$ ?

Set  $\sigma^2 = \text{Var}(X)$ . Since the  $X_i$  are iid copies of r.v.  $X$ ,

$$E[(\bar{X}_S - \mu)^2] = \text{Var}(\bar{X}_S) = \frac{1}{S^2} \sum_{i=1}^S \text{Var}(X_i) = \frac{1}{S} \sigma^2$$

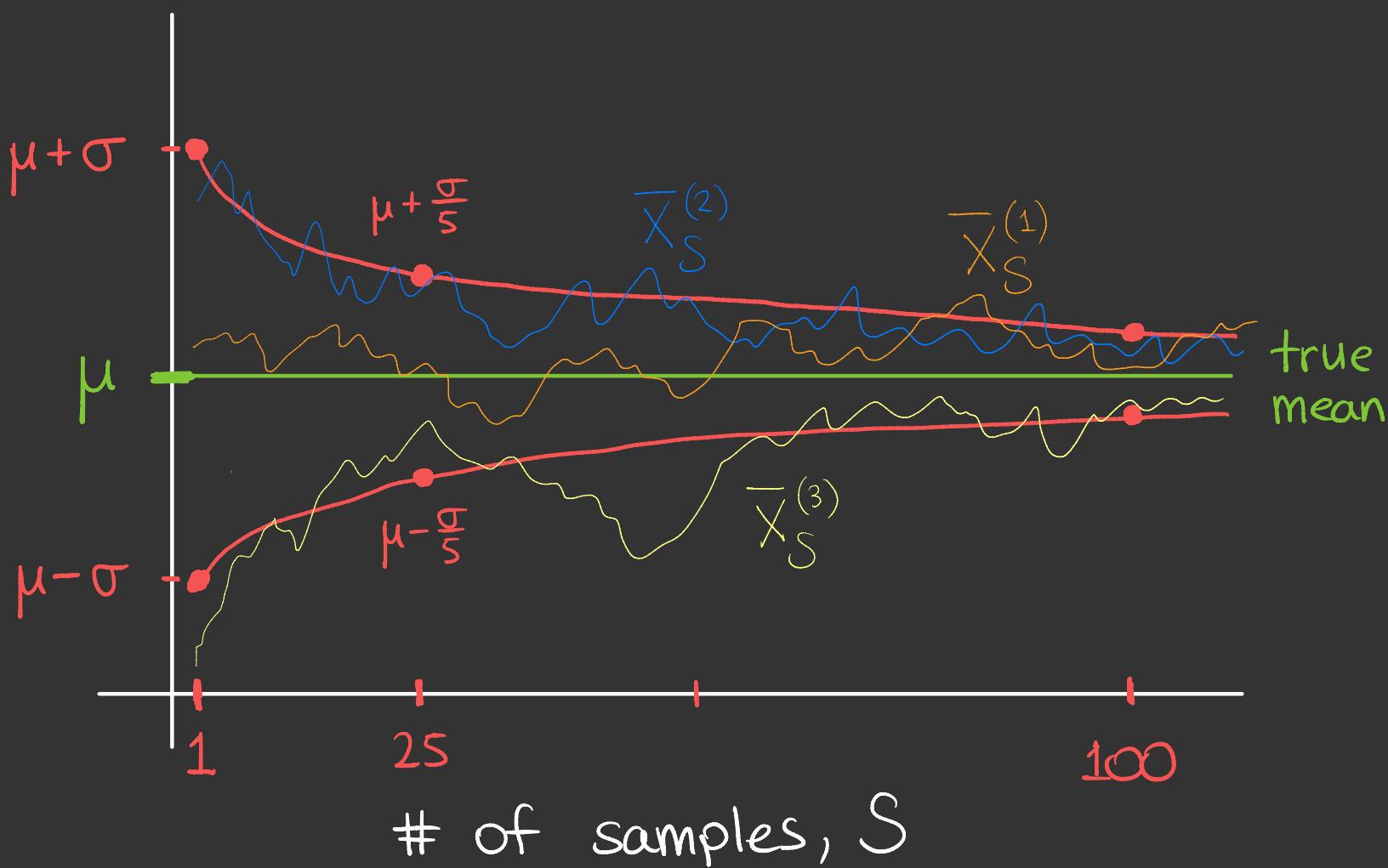
$$\implies \text{Root mean squared error} = \frac{\sigma}{\sqrt{S}}.$$

- If we approximate  $\sigma^2 = \text{Var}(X)$  by

$$\hat{\sigma}_S^2 = \frac{1}{S-1} \sum_{i=1}^S (x_i - \bar{x}_S)^2$$

then the Monte Carlo standard error is

$$\frac{\hat{\sigma}_S}{\sqrt{S}}.$$



- The Central Limit Theorem (CLT) gives an even more complete characterization of MC error.

CLT. If  $X_1, X_2, \dots$  iid copies of  $X$  with  $E[X] = \mu$  and  $\text{Var}(X) = \sigma^2 < \infty$ , then

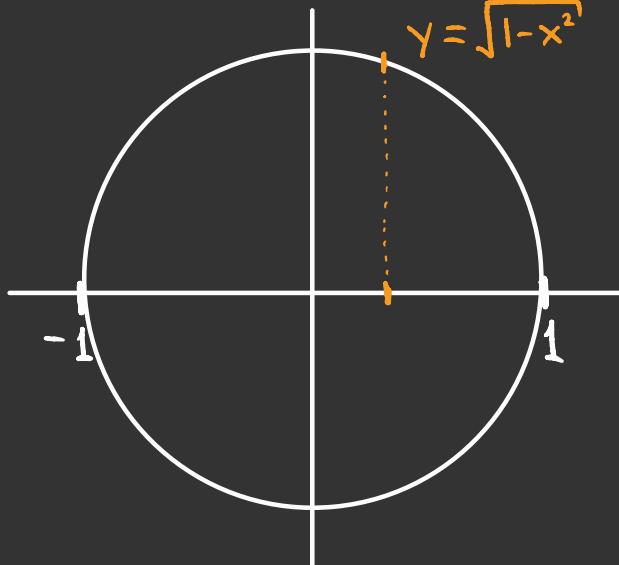
$$\sqrt{S}(\bar{X}_S^* - \mu) \xrightarrow{\text{distribution}} N(0, \sigma^2).$$

Equivalently,  $\frac{\sqrt{S}}{\sigma}(\bar{X}_S - \mu) \xrightarrow{d} N(0, 1)$ .

\* Here,  $\bar{X}_S = \frac{1}{S} \sum_{i=1}^S X_i$

## A non-Bayesian example.

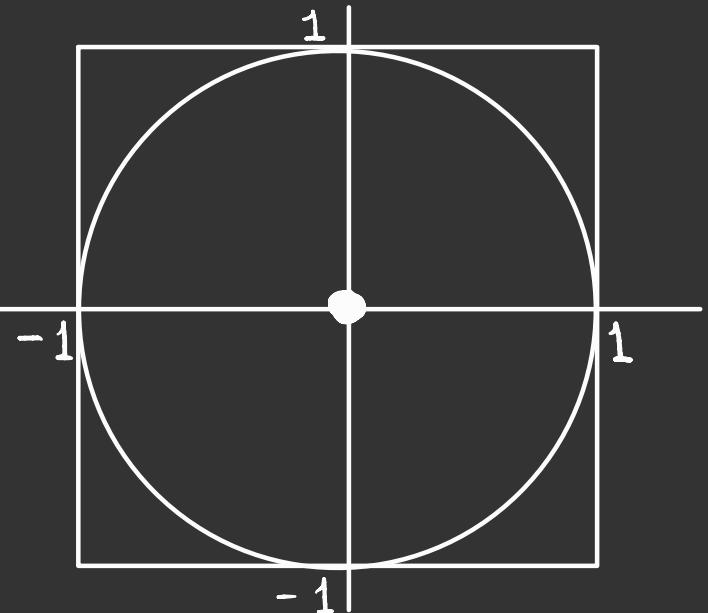
- Goal : Find the area of a circle of radius 1.
- One approach is to use calculus



$$\text{Area} = 4 \int_0^1 \sqrt{1 - x^2} dx$$

Use trig rules to solve.

- Monte Carlo approach



Circle or ball of radius 1

- Let  $B = \{(x, y) : x^2 + y^2 < 1\}$
- Let  $S$  be the square around  $B$  (side length 2)
- Let  $X, Y \stackrel{\text{iid}}{\sim} \text{Uniform}(-1, 1)$
- Then :

$$\begin{aligned}\text{Area}(B) &= P((X, Y) \in B) \text{Area}(S) \\ &= 4 P((X, Y) \in B).\end{aligned}$$

$$\text{Area}(B) = 4 \mathbb{P}((X, Y) \in B) = 4 \mathbb{E}[1_{\{(X, Y) \in B\}}]$$

By LLN, if we draw  $S$  samples  $X_i \stackrel{\text{iid}}{\sim} \text{Unif}(-1, 1)$  and  $Y_i \stackrel{\text{iid}}{\sim} \text{Unif}(-1, 1)$ , then for large  $S$

$$\begin{aligned}\text{Area}(B) &\approx \frac{1}{S} \sum_{i=1}^S 4 \mathbb{1}_{\{(X_i, Y_i) \in B\}} \\ &= \frac{1}{S} \sum_{i=1}^S 4 \mathbb{1}_{(X_i^2 + Y_i^2 < 1)}.\end{aligned}$$

More generally, to approximate volume of unit ball

$$B_1(0) = \left\{ \vec{x} \in \mathbb{R}^d : \sum_{i=1}^d x_i^2 < 1 \right\},$$

generate  $S \times d$  matrix

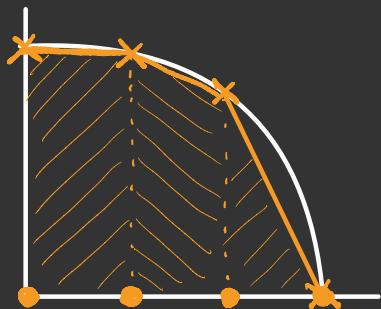
$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & & & \vdots \\ x_{S1} & x_{S2} & \cdots & x_{Sd} \end{pmatrix}$$

where each  $x_{ij} \stackrel{\text{iid}}{\sim} \text{Unif}(-1, 1)$ .

Then

$$\text{Vol}(B_1(o)) \approx \frac{2^d}{S} \sum_{l=1}^S \mathbb{1}\left(X_{l1}^2 + \dots + X_{ld}^2 < 1\right).$$

# Monte Carlo vs. deterministic grid



- Grid approach to numerical integration

$$\int_a^b f(x) dx$$

- Split  $[a,b]$  into  $(x_0, x_1, \dots, x_g)$   
where  $x_0 = a$  and  $x_g = b$  and  
 $x_i < x_{i+1} \quad \forall i = 0, \dots, g-1$ .

- Then

$$\int_a^b f(x) dx \approx \sum_{i=1}^{g-1} f(x_i)(x_{i+1} - x_i)$$

In  $d$ -dimensions, have  $g^d$  grid points!

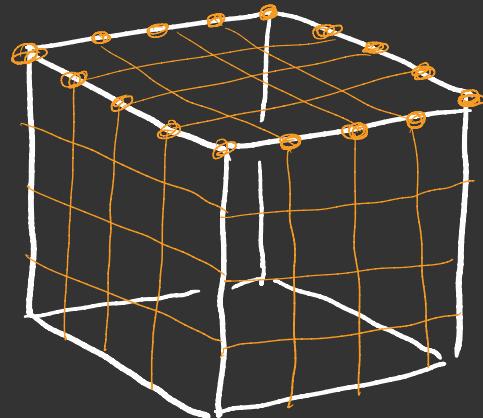
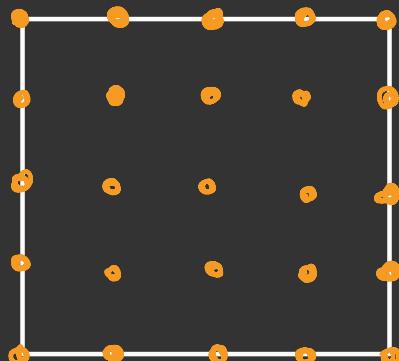
$$5^1 = 5$$

$$5^2 = 25$$

$$5^3 = 125$$

$$5^{10} \approx 10,000,000$$

$$5^{20} \approx 100 \text{ trillion}$$



$$\int_a^b f(x) dx$$

$$\int_{a_1}^{b_1} \int_{a_2}^{b_2} f(x, y) dx dy$$

$$\int_{a_1}^{b_1} \int_{a_2}^{b_2} \int_{a_3}^{b_3} f(x, y, z) dx dy dz$$

Monte Carlo convergence rate is  $\frac{1}{\sqrt{S}}$ .

Depends only on # of samples, "dimension free"

## Beta-binomial example.

- Prior:  $p(\theta) = \text{Beta}(1, 1)$
- Observations  $y_1, \dots, y_{10} | \theta \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$  with  
 $\sum y_i = 3$
- $p(\theta | \vec{y}_n) = \text{Beta}(4, 8)$
- Since the posterior has closed form, we can compute certain posterior quantities exactly. Let's compare these to Monte

Carlo approximations.

- Posterior mean :  $E[\theta | \vec{y}_n] = \frac{a + \sum y_i}{a + b + n}$
- Posterior variance
- $P(\theta < \frac{3}{4} | y)$
- A 0.1 quantile

## Posterior inference of arbitrary functions

If interested in a function  $g(\theta)$  rather than  $\theta$  itself, can generate iid samples via

for  $1 \leq i \leq S$ :

- Sample  $\theta_i \stackrel{\text{iid}}{\sim} p(\theta | y)$
- Set  $y_i = g(\theta_i)$ .

Then  $y_1, \dots, y_S \stackrel{\text{iid}}{\sim} p(g(\theta) | y)$ .

- Example (see Hoff, Section 4.2 for details) :

$$\text{log odds } (\theta) = \log \left( \frac{\theta}{1-\theta} \right)$$

## Sampling from joint posterior

If  $\vec{\Theta}_d = (\Theta_1, \dots, \Theta_d) \in \mathbb{R}^d$  AND the parameters  $\Theta_j$  are conditionally independent given the data  $y$ , then we can generate iid samples  $\vec{\Theta}_{d,i} \stackrel{\text{iid}}{\sim} p(\vec{\Theta}_d | y)$  via

- for  $1 \leq i \leq S$  :
  - for  $1 \leq j \leq d$  :
    - draw  $\Theta_{ij} \sim p(\Theta_j | y)$ .

## Sampling from posterior predictive distribution

PPD is (assuming  $Y_1, \dots, Y_{n+1}$  cond. ind. given  $\Theta$ )

$$p(Y_{n+1} | \vec{y}_n) = \int_{\Theta} p(Y_{n+1} | \theta) p(\theta | \vec{y}_n) d\theta.$$

for  $1 \leq i \leq S$ :

- Draw  $\theta_i \sim p(\theta | \vec{y}_n)$
- Draw  $\tilde{y}_i \sim p(Y_{n+1} | \theta_i)$

$\Rightarrow \{(\theta_1, \tilde{y}_1), \dots, (\theta_S, \tilde{y}_S)\} = S$  independent

samples from  $p(\theta, y_{n+1} | \vec{y}_n)$  and

$\{\tilde{y}_1, \dots, \tilde{y}_s\}$  iid samples from  $p(y_{n+1} | \vec{y}_n)$ .

## Posterior predictive model checking

- How plausible is our model?
- Suppose we're interested in some test statistic  $t(\vec{Y}_n)$  of the data, e.g.,

$$t(\vec{Y}_n) = \frac{1}{n} \sum_{i=1}^n Y_i.$$

If our proposed Bayesian model

$$p(\theta | \vec{y}_n) \propto p(\vec{y}_n | \theta) p(\theta)$$

is reasonable, the test values  $t(\vec{Y}_n^{ppd})$   
w/ datasets  $\vec{Y}_n^{ppd}$  drawn from the PPD

$$p(\vec{Y}_n^{ppd} | \vec{y}_n) = \int_{\Theta} p(\vec{Y}_n^{ppd} | \theta) p(\theta | \vec{y}_n) d\theta$$

Should resemble the observed  $t(\vec{Y}_n)$ .

## Basic idea:

- Simulate datasets from PPD.
- Evaluate  $t$  on simulated datasets.
- Check whether observed  $t$  matches those evaluated on simulated data.

## Monte Carlo approach.

for  $1 \leq i \leq S$ :

- Sample  $\Theta_i \sim p(\theta | \vec{y}_n)$
- Sample  $\vec{y}_{n,i}^{\text{ppd}} = (y_{1,i}^{\text{ppd}}, \dots, y_{n,i}^{\text{ppd}}) \stackrel{\text{iid}}{\sim} p(y | \Theta_i)$
- Compute  $t_i = t(\vec{y}_{n,i}^{\text{ppd}})$

Return: Samples  $\{t_1, \dots, t_S\}$  from PPD of  $t(\vec{Y}_n^{\text{ppd}})$

See dental health code for example.

End Ch. 4