

Context

Variance reduction techniques are used to improve the efficiency of Monte Carlo methods for approximating integrals. In this work we propose an adaptive stratified sampling approach based on theoretical bounds of the strata variances using the monotonicity property of the function of interest.

Problem statement

Consider the problem of estimating

$$\bar{\phi} = \mathbb{E}(\phi(X)) \quad (1)$$

where

- $\phi : \mathbb{R}^+ \rightarrow [0, 1]$
- ϕ is continuous
- ϕ is monotonically increasing
- $\phi(0) = 0$ and $\lim_{x \rightarrow \infty} \phi(x) = 1$
- The CDF F_X of X is known

Motivation

- In the field of microbiology, Quantitative Risk Assessment (QRA) models are used for estimating the risk of a food borne disease
- QRA models include monotonically increasing functions w.r.t the initial concentration of bacteria
- Given the cost of function evaluation the aim is to reduce the sampling budget compared to the simple Monte Carlo algorithm

How stratification works

Consider the N strata $S_i = [l_i, u_i]$ for $i = 1, 2, \dots, N$ with:

- stratum probability $\omega_i = P[X \in S_i]$
- stratum variance $\tau_i^2 = \text{Var}(\phi(X)|X \in S_i)$

Then the stratified sampling estimator can be written as

$$\hat{\phi}^{\text{ST}} = \sum_{i=1}^N \frac{\omega_i}{n_i} \sum_{j=1}^{n_i} \phi(X_{i,j}) \quad (2)$$

and the variance of this estimator

$$\text{Var}(\hat{\phi}^{\text{ST}}) = \sum_{i=1}^N \frac{\omega_i^2 \tau_i^2}{n_i} \quad (3)$$

with

- stratum sample size n_i
- stratum samples $X_{i,j}, j = 1, 2, \dots, n_i$

The optimal choice for n_i , with total budget $\sum_{i=1}^N n_i = n$, obtained by minimizing the variance is

$$n_i = \frac{\omega_i \tau_i}{\sum_{i=1}^N \omega_i \tau_i} n \quad (4)$$

The variance of the stratified sampling estimator with optimal allocation of sampling budget is

$$\text{Var}(\hat{\phi}_{\text{opt}}^{\text{ST}}) = \frac{1}{n} \left(\sum_{i=1}^N \omega_i \tau_i \right)^2 \quad (5)$$

Stratification vs Simple Monte Carlo

- The simple Monte Carlo estimate is a special case of the stratified sampling estimate with $N = 1$
- Both estimates are unbiased
- The variance of the Monte Carlo estimator can be shown to be larger than $\text{Var}(\hat{\phi}_{\text{opt}}^{\text{ST}})$

Conservative approach

- The stratum variances τ_i^2 are unknown
- We propose using upper bounds instead of pilot sample estimates

Popoviciu's inequality provides an upper bound on the variance of bounded random variables

$$\tau_{n,i}^2 \leq \frac{1}{4}(\phi(u_i) - \phi(l_i))^2 = \frac{1}{4}\Delta_i^2 \quad (6)$$

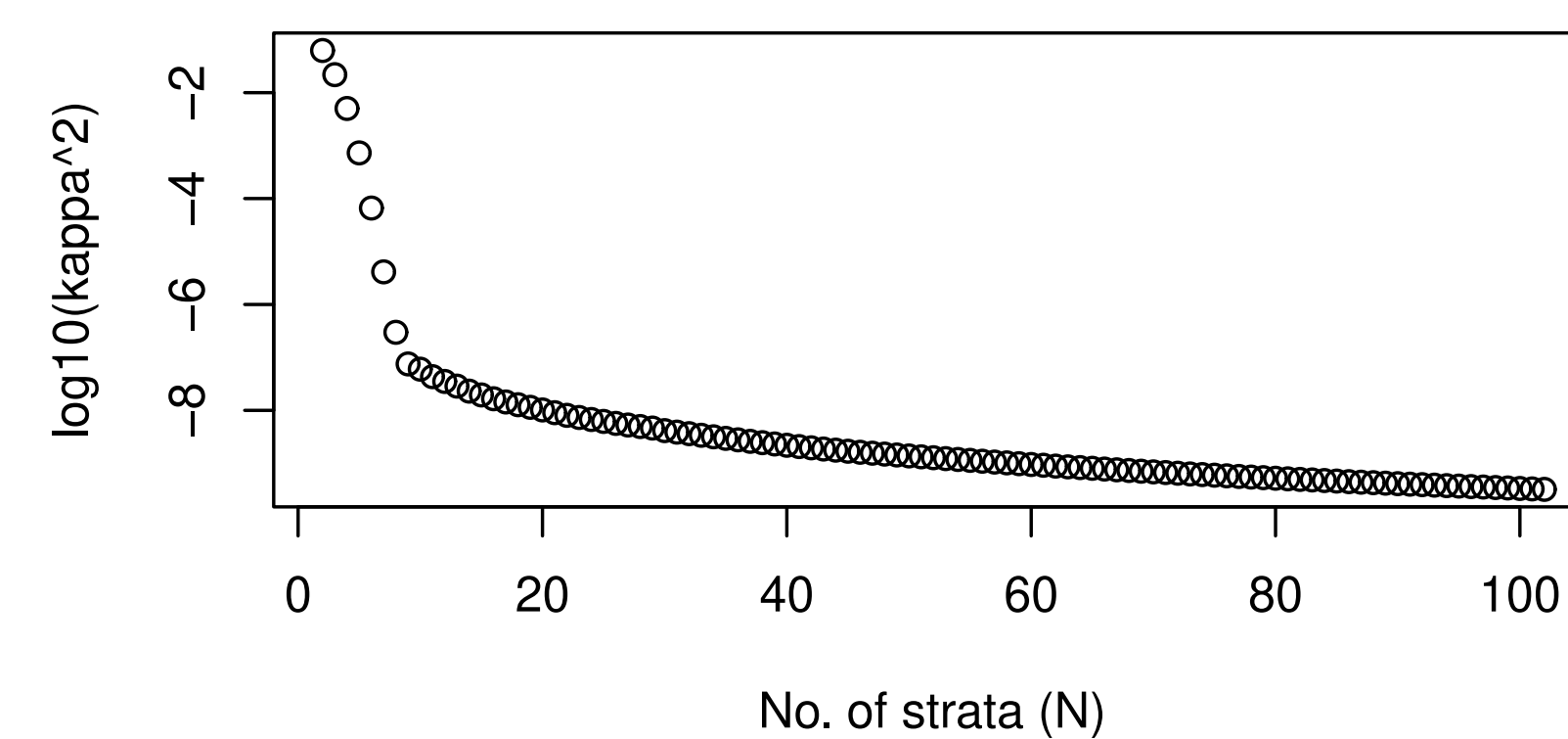
- An upper bound of the optimal variance can be written as

$$\text{Var}(\hat{\phi}_{\text{opt}}^{\text{ST}}) \leq \frac{1}{4n} \left(\sum_{i=1}^N \omega_i \Delta_i \right)^2 = \frac{\kappa^2}{n} \quad (7)$$

- For fixed n , the upper bound $\frac{\kappa^2}{n}$ decreases as the number of strata N increases
- It can be shown that a split in the i -th stratum with any split proportion $0 < \alpha = \frac{\omega_{1,i}}{\omega_i} < 1$ and $0 < \beta = \frac{\Delta_{1,i}}{\Delta_i} < 1$, reduces κ

$$\omega_i \Delta_i > \omega_{1,i} \Delta_{1,i} + (\omega_i - \omega_{1,i})(\Delta_i - \Delta_{1,i}) \quad (8)$$

Upper bound on variance



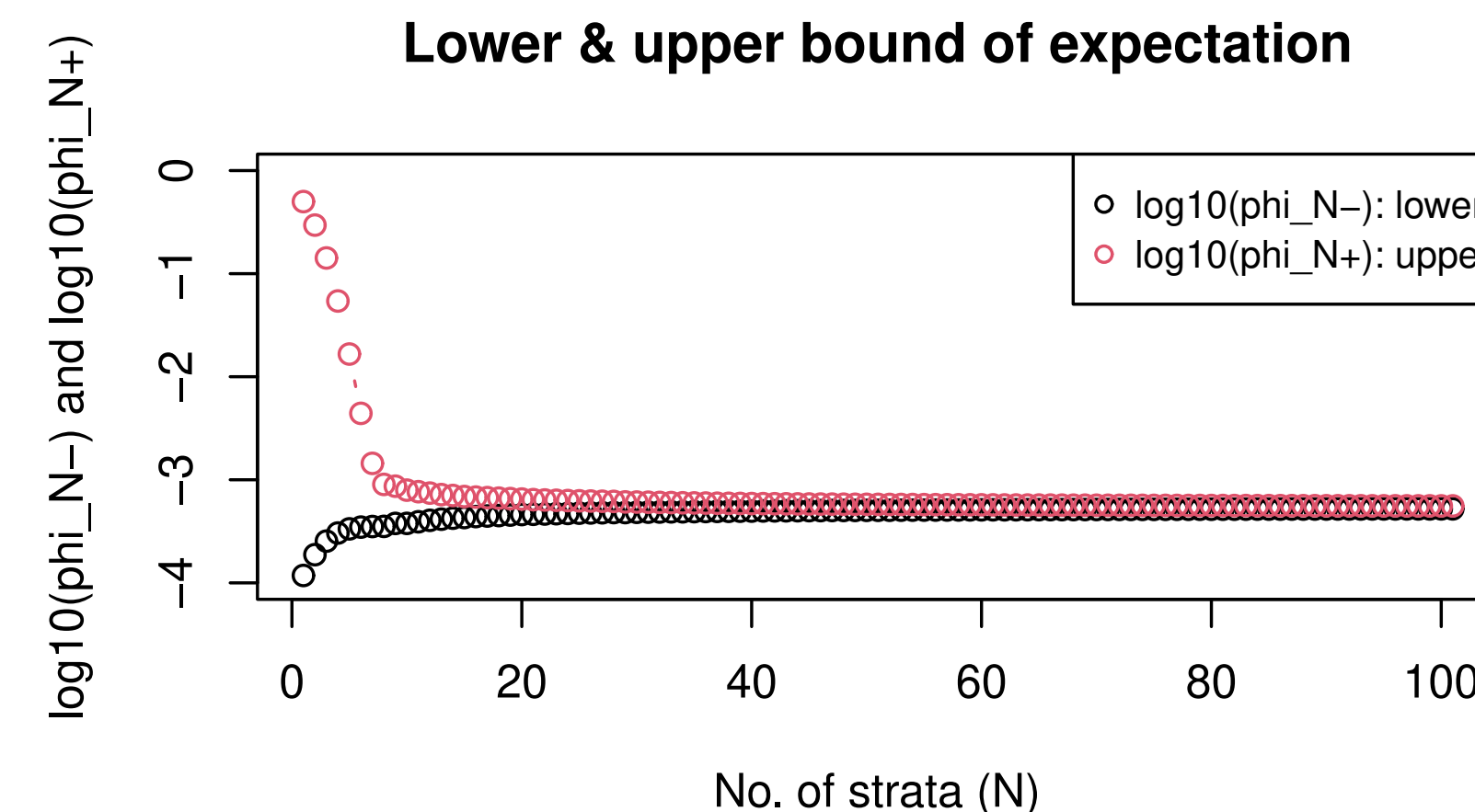
- Bounds on the expectation can be derived using the monotonicity property of ϕ

$$\bar{\phi} \geq \sum_{i=1}^N \omega_i \phi(l_i) = \phi_N^- \quad (9)$$

$$\bar{\phi} \leq \sum_{i=1}^N \omega_i \phi(u_i) = \phi_N^+$$

- The upper and lower bounds for the expectation converges as number of strata N increases

Lower & upper bound of expectation



Adaptive algorithm

- The idea is to split the stratum with maximum $\omega_i \Delta_i$ value (the highest contribution to the variance upper bound)
- A split which divides either ω_i or Δ_i into half, results in reducing the contribution of that particular stratum also by half
- To obtain a split with α or β close to 0.5 we propose the following strategy

$$\begin{aligned} k &\leftarrow \arg \max_i (\omega_i \Delta_i) \\ \text{if } k = N &\text{ then} \\ &X_{(k)} = 2X_{(k-1)} \\ \text{else} \\ &X_{(k)} = \frac{X_{(k-1)} + X_{(k)}}{2} \end{aligned}$$

- For a fixed number of strata N the **sampling budget** n can be obtained by fixing the upper bound of coefficient of variation by δ

$$\begin{aligned} \text{CV}(\hat{\phi}^{\text{ST}}) &= \frac{\sqrt{\text{Var}(\hat{\phi}^{\text{ST}})}}{\bar{\phi}} \leq \frac{\kappa}{\sqrt{n} \bar{\phi}_N} \leq \delta \\ \Rightarrow n &= \left\lceil \frac{\kappa^2}{(\bar{\phi}_N)^2 \delta^2} \right\rceil \end{aligned} \quad (10)$$

- The optimal strata sizes given by (4) (by substituting Δ_i) might result in non integer values
- To ensure each stratum has at least one sample, we recompute the budget as n^{corr} by taking the ceiling value

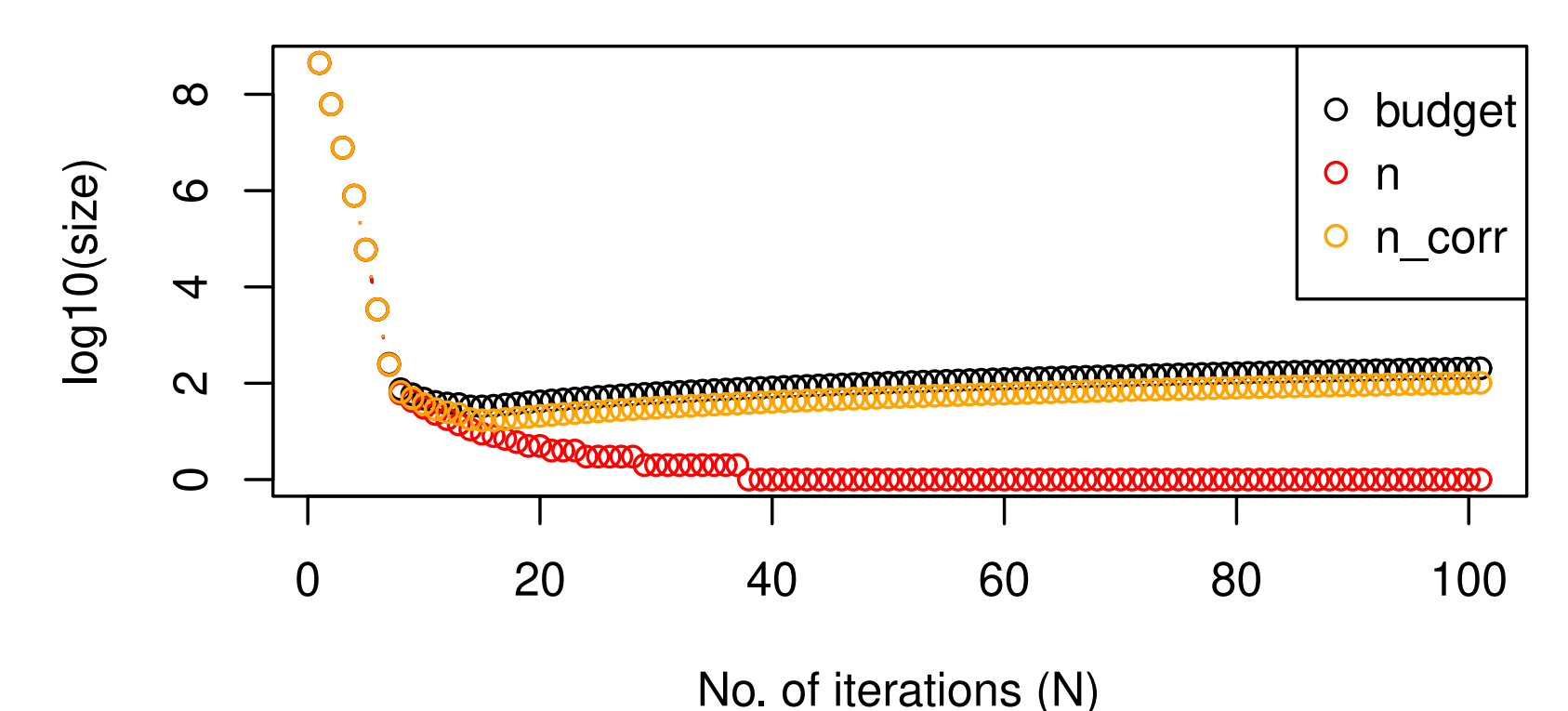
$$n^{\text{corr}} = \sum_{i=1}^N \tilde{n}_i = \sum_{i=1}^N \lceil n_i \rceil \geq n \quad (11)$$

- The upper bound for variance with the corrected sample size is smaller

$$\sum_{i=1}^N \frac{\omega_i^2 \Delta_i^2}{\tilde{n}_i} \geq \sum_{i=1}^N \frac{\omega_i^2 \Delta_i^2}{n_i} \quad (12)$$

- The actual **budget** of the algorithm is n^{corr} + the additional evaluations made in each strata for splitting

Sample sizes and budget



Stopping rule: The algorithm stops splitting the strata as the budget starts increasing

Experimental results

- The proposed method is implemented and compared to simple Monte Carlo for estimating the risk in a QRA model
- A simple Monte Carlo approach requires a budget of 2503 samples to achieve a 10% CV
- The proposed stratified sampling approach requires 33 samples to achieve the same CV