

# Monte Carlo Methods

Introduction

Generating Samples

Markov Chain Monte Carlo + Discrepancy

Improving Efficiency

Selected Topics

Git website and repository  
Canvas

Fred Hickernell, Fall 2025

Updated 2025 December 1



# Selected Topics

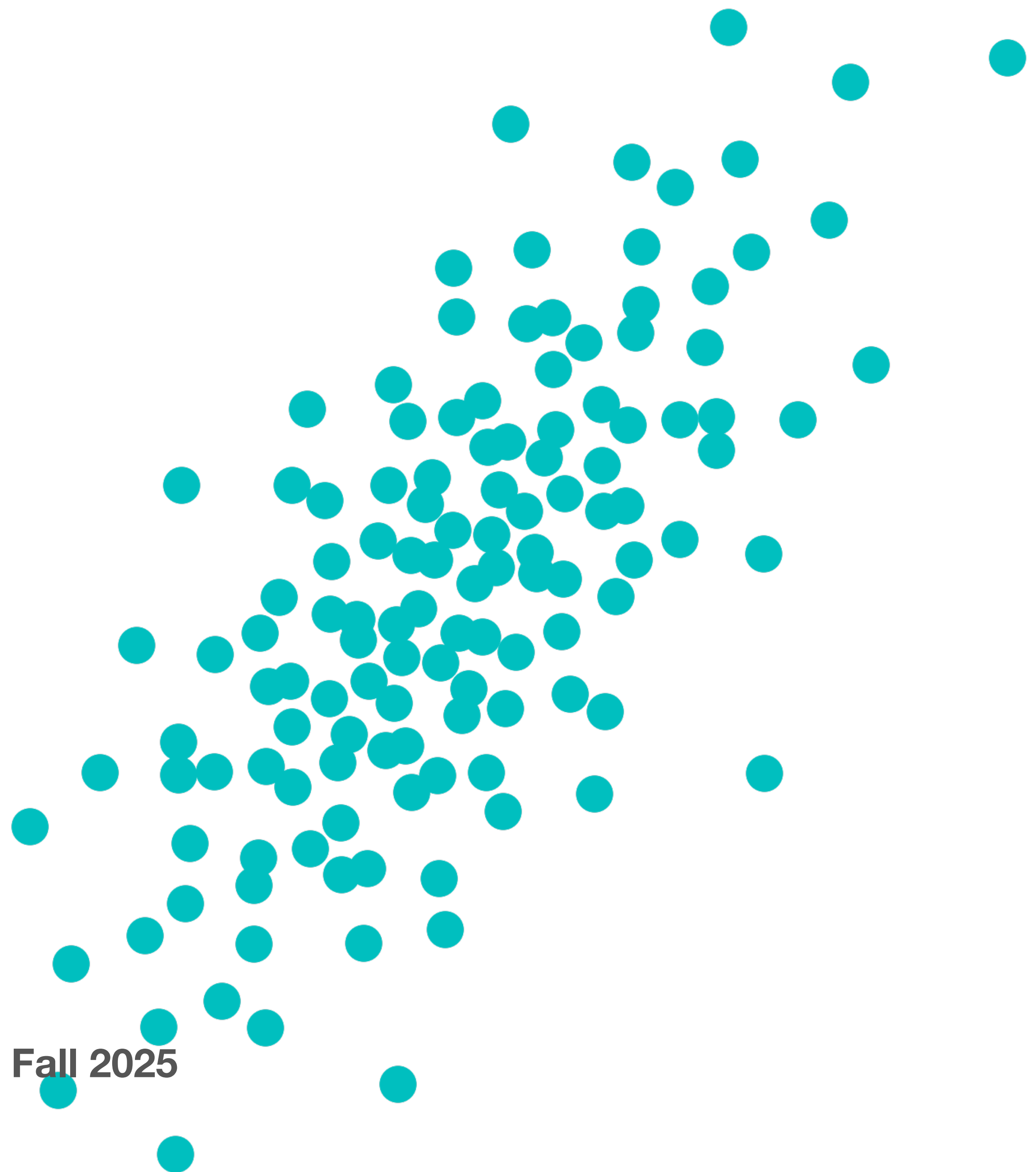
Owen, Chapters ???

Project Presentations Nov 24–25

Take-Home Final Available Dec 3

Take-Home Final Due and  
In-Class Final Dec 11, 2–4PM,  
WH 115

**MATH 565 Monte Carlo Methods, Fred Hickernell, Fall 2025**





# Parallel Computing

Many personal computers now have GPUs built in. This allows one to perform computations in parallel using PyTorch. Here are some things to keep in mind

- Monte Carlo calculations are **pleasantly parallel** because the function values can be obtained independently, and one can even compute multiple expectations independently.
- Calculations on GPUs are typically using **32-bit floating point**, not 64-bit floating point. This means that we have about 6 significant digits of accuracy, not 15.
- For small problems, GPUs may not be better than CPUs, but for **large numbers of samples** or **large numbers of similar problems**, the speed-up can be impressive.



# Gradient Descent

Optimization problems take the form

$$\boldsymbol{\theta}_* = \underset{\boldsymbol{\theta} \in \mathbb{R}^m}{\operatorname{argmin}} \operatorname{Loss}(\boldsymbol{\theta})$$

Gradient descent updates the estimate of  $\boldsymbol{\theta}_*$  iteratively by moving **downhill**:

$$\boldsymbol{\theta}_0 \text{ given, } \boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \nabla \operatorname{Loss}(\boldsymbol{\theta}_k),$$

where  $\eta$  is the step size or learning rate.

- This will often converge since  $\nabla \operatorname{Loss}(\boldsymbol{\theta}_k) \rightarrow \mathbf{0}$  as  $\boldsymbol{\theta}_k \rightarrow \boldsymbol{\theta}_*$
- But it may converge to a **local minimum**
- This may require more steps than methods that use **second derivatives**, but each step should be cheaper and easier to compute



# Stochastic Gradient Descent

Loss function in **regression** often take the form

$$\text{Loss}(\boldsymbol{\theta}) = \sum_{i=1}^N \text{loss}_i(\boldsymbol{\theta}; \text{data}), \quad \text{e.g., } \text{loss}_i(\boldsymbol{\theta}; \text{data}) = [y_i - \theta_1 \exp(\theta_2 x_i)]^2$$

The loss to be optimized is a **sum** of many losses. If  $N$  is large, then the gradient is **expensive**. An alternative is **stochastic gradient descent**

$$\boldsymbol{\theta}_0 \text{ given, } \boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \nabla \text{loss}_{i_k}(\boldsymbol{\theta}_k)$$

where  $i_0, i_1, \dots$  are chosen IID uniform over  $1, \dots, N$ , which is **cheap**.





# Stochastic Gradient Descent

Loss function in **regression** often take the form

$$\text{Loss}(\boldsymbol{\theta}) = \sum_{i=1}^N \text{loss}_i(\boldsymbol{\theta}; \text{data}), \quad \text{e.g., } \text{loss}_i(\boldsymbol{\theta}; \text{data}) = [y_i - \theta_1 \exp(\theta_2 x_i)]^2$$

gradient is **expensive**. An alternative is **stochastic gradient descent**

**vanilla**  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \nabla \text{loss}_{i_k}(\boldsymbol{\theta}_k)$

**decaying  $\eta$**   $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta_k \nabla \text{loss}_{i_k}(\boldsymbol{\theta}_k), \quad \eta_k = \frac{\eta_0}{\sqrt{1 + k/100}}$

**mini-batch**  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \frac{1}{B} \sum_{i_{kj}=1}^B \nabla \text{loss}_{i_{kj}}(\boldsymbol{\theta}_k), \quad i_{kj} \stackrel{\text{iid}}{\sim} \mathcal{U}\{1, \dots, N\}$

**mini-batch + decay**  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta_k \frac{1}{B} \sum_{i_{kj}=1}^B \nabla \text{loss}_{i_{kj}}(\boldsymbol{\theta}_k), \quad i_{kj} \stackrel{\text{iid}}{\sim} \mathcal{U}\{1, \dots, N\} \quad \eta_k = \frac{\eta_0}{\sqrt{1 + k/100}}$



# Two Level Monte Carlo

Let's start with a simple example:  $Y = Y_1 + Y_2$  and we want to estimate

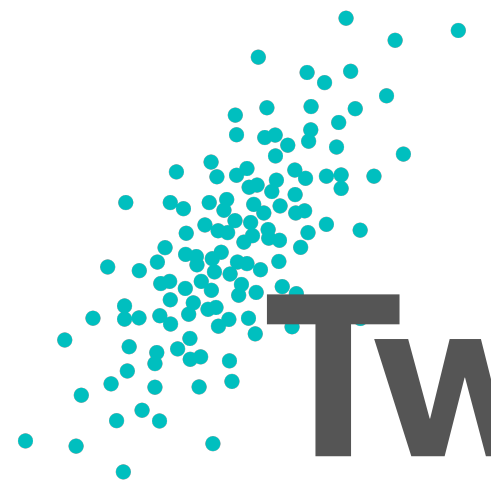
$$\mu := \mu_1 + \mu_2, \text{ where } \mu := \mathbb{E}(Y), \mu_1 := \mathbb{E}(Y_1), \mu_2 := \mathbb{E}(Y_2)$$

using IID samples

$$\hat{\mu} := \hat{\mu}_1 + \hat{\mu}_2, \quad \hat{\mu}_1 := \frac{1}{n_1} \sum_{i=0}^{n_1-1} Y_{i1}, \quad \hat{\mu}_2 := \frac{1}{n_2} \sum_{i=0}^{n_2-1} Y_{i2}, \quad Y_{0,1}, Y_{1,1}, \dots, Y_{0,2}, Y_{1,2}, \text{ IID}$$

$$\text{mse}(\hat{\mu}) = \frac{\text{var}(Y_1)}{n_1} + \frac{\text{var}(Y_2)}{n_2}$$

and the cost of one  $Y_\ell$  is  $\$_\ell$ . How does one choose  $n_1$  and  $n_2$  optimally to minimize  $\text{mse}(\hat{\mu})$ ?



# Two Level Monte Carlo

$$Y = Y_1 + Y_2, \mu := \mu_1 + \mu_2, \mu := \mathbb{E}(Y), \mu_1 := \mathbb{E}(Y_1), \mu_2 := \mathbb{E}(Y_2)$$

$$\hat{\mu} := \hat{\mu}_1 + \hat{\mu}_2, \quad \hat{\mu}_1 := \frac{1}{n_1} \sum_{i=0}^{n_1-1} Y_{i1}, \quad \hat{\mu}_2 := \frac{1}{n_2} \sum_{i=0}^{n_2-1} Y_{i2}, \quad Y_{0,1}, Y_{1,1}, \dots, Y_{0,2}, Y_{1,2}, \text{ IID}$$

$$\begin{aligned} \text{mse}(\hat{\mu}) &= \frac{\text{var}(Y_1)}{n_1} + \frac{\text{var}(Y_2)}{n_2} \\ &= \underbrace{n_1 \$1}_{\text{time spent on } Y_1} \underbrace{\frac{\text{var}(Y_1)}{n_1^2 \$1}}_{\text{mse}(\hat{\mu}_1) \text{ per time spent on } Y_1} + \underbrace{n_2 \$2}_{\text{time spent on } Y_2} \underbrace{\frac{\text{var}(Y_2)}{n_2^2 \$2}}_{\text{mse}(\hat{\mu}_2) \text{ per time spent on } Y_2} \end{aligned}$$

To optimize make  $\frac{\text{var}(Y_1)}{n_1^2 \$1} = \frac{\text{var}(Y_2)}{n_2^2 \$2}$





# Multilevel Monte Carlo

$$Y = Y_1 + \dots + Y_L, \mu := \mu_1 + \dots + \mu_L, \mu := \mathbb{E}(Y), \mu_\ell := \mathbb{E}(Y_\ell)$$

$$\hat{\mu} := \hat{\mu}_1 + \dots + \hat{\mu}_L \quad \hat{\mu}_\ell := \frac{1}{n_\ell} \sum_{i=0}^{n_\ell-1} Y_{i1}, \quad Y_{0,1}, Y_{1,1}, \dots, Y_{0,2}, Y_{1,2}, \dots, \dots \text{ IID}$$

$$\text{mse}(\hat{\mu}) = \frac{\text{var}(Y_1)}{n_1} + \dots + \frac{\text{var}(Y_L)}{n_L}, \quad \$_\ell = \text{cost of one } Y_\ell$$

$$= \underbrace{n_1 \$1}_{\text{time spent on } Y_1} \underbrace{\frac{\text{var}(Y_1)}{n_1^2 \$1}}_{\text{mse}(\hat{\mu}_1) \text{ per time spent on } Y_1} + \dots + \underbrace{n_L \$L}_{\text{time spent on } Y_L} \underbrace{\frac{\text{var}(Y_L)}{n_L^2 \$L}}_{\text{mse}(\hat{\mu}_L) \text{ per time spent on } Y_L}$$

$$\text{To optimize make } \frac{\text{var}(Y_1)}{n_1^2 \$1} = \dots = \frac{\text{var}(Y_L)}{n_L^2 \$L}, \text{ i.e. } n_\ell \propto \sqrt{\frac{\text{var}(Y_\ell)}{\$_\ell}}$$



# Multilevel Monte Carlo

- Larger  $\text{var}(Y_\ell)$  implies **larger**  $n_\ell$
- Larger  $\$_\ell$  implies **smaller**  $n_\ell$
- $\text{var}(Y_\ell)$  and  $\$_\ell$  may need to be estimated by **pilot samples**
- Often  $\mu := \mu_1 + \cdots + \mu_L + \Delta$ , where  **$\Delta$**  is small but cannot be sampled
- For the optimal sample sizes  $n_\ell = c\sqrt{\text{var}(Y_\ell)/\$_\ell}$

$$\text{total cost}(\hat{\mu}) = c \left[ \sqrt{\text{var}(Y_1) \$_1} + \cdots + \sqrt{\text{var}(Y_L) \$_L} \right]$$

$$\text{mse}(\hat{\mu}) = \frac{1}{\text{total cost}(\hat{\mu})} \left[ \sqrt{\text{var}(Y_1) \$_1} + \cdots + \sqrt{\text{var}(Y_L) \$_L} \right]^2$$



# Multilevel Monte Carlo for $\infty$ -D Expectations

$\mu = \lim_{d \rightarrow \infty} \mathbb{E}[f_d(\mathbf{X}_{1:d})]$ , e.g.,  $f_d$  is an option payoff using  $d$  time steps

Let  $Y_\ell = f_{d_\ell}(\mathbf{X}) - f_{d_{\ell-1}}(\mathbf{X}_{1:d_{\ell-1}})$ , for  $\mathbf{X} \sim \mathcal{U}[0,1]^{d_\ell}$ , where  $f_{d_0} = 0$ . Then

$$\begin{aligned} \mu = \mathbb{E} & \left[ \underbrace{f_{d_1}(\mathbf{X}_{1:d_1}) - f_{d_0}}_{Y_1} + \underbrace{f_{d_2}(\mathbf{X}_{1:d_2}) - f_{d_1}(\mathbf{X}_{1:d_1})}_{Y_2} + \cdots \right. \\ & \left. + \underbrace{f_{d_L}(\mathbf{X}_{1:d_L}) - f_{d_{L-1}}(\mathbf{X}_{1:d_{L-1}})}_{Y_L} \right] + \lim_{d \rightarrow \infty} \mathbb{E}[f_d(\mathbf{X}_{1:d}) - f_{d_L}(\mathbf{X}_{1:d_L})] \end{aligned}$$

$$\hat{\mu} := \hat{\mu}_1 + \cdots + \hat{\mu}_L, \quad \hat{\mu}_\ell := \frac{1}{n_\ell} \sum_{i=0}^{n_\ell-1} Y_{i\ell}, \quad Y_{0,1}, Y_{1,1}, \dots, Y_{0,2}, Y_{1,2}, \text{ IID}$$

Typically  $n_\ell \propto d_\ell$ . If  $\text{var}(Y_{d_\ell})$  decays, then MLMC works well by choosing  $d_1 < d_2 < \cdots < d_L$  with  $d_L$  large enough.



# Multilevel Monte Carlo for $\infty$ -D Expectations

$$\mu = \lim_{d \rightarrow \infty} \mathbb{E}[f_d(X_{1:d})], \quad Y_\ell = f_{d_\ell}(X) - f_{d_{\ell-1}}(X_{1:d_{\ell-1}}), \text{ for } X \sim \mathcal{U}[0,1]^{d_\ell}, f_{d_0} = 0.$$

$$\hat{\mu} := \hat{\mu}_1 + \cdots + \hat{\mu}_L, \quad \hat{\mu}_\ell := \frac{1}{n_\ell} \sum_{i=0}^{n_\ell-1} Y_{i\ell}, \quad Y_{0,1}, Y_{1,1}, \dots, Y_{0,2}, Y_{1,2}, \text{ IID}$$

If  $d_\ell = m^\ell$ ,  $\text{var}(Y_\ell) \leq \beta v^\ell$ , and  $\$_\ell \leq \gamma d_\ell = \gamma m^\ell$ , then

$$\begin{aligned} \text{total cost}(\hat{\mu}) &\leq c \left[ \sqrt{\beta \gamma v m} + \cdots + \sqrt{\beta \gamma (v m)^L} \right] \\ &= c \sqrt{\beta \gamma} \left[ (v m)^{1/2} + \cdots + (v m)^{L/2} \right] \\ &\leq \begin{cases} \frac{c \sqrt{\beta \gamma v m}}{1 - \sqrt{v m}} \propto \text{cost of } Y_1 \text{ part}, & v m < 1 \\ \frac{c \sqrt{\beta \gamma (v m)^L}}{1 - 1/\sqrt{v m}} \propto \text{cost of } Y_L \text{ part}, & v m > 1 \end{cases} \end{aligned}$$