



# UCL

## Nested expectations with kernel quadrature

Dr François-Xavier Briol  
Department of Statistical Science  
University College London

<https://fxbriol.github.io/>



Hudson Chen



Masha Naslidnyk

# Nested expectations

Quantity of interest:

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta.$$

# Nested expectations

Quantity of interest:

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta.$$

inner expectation  
(against  $\mathbb{P}_{\theta}$ )

# Nested expectations

Quantity of interest:

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta .$$

inner expectation  
(against  $\mathbb{P}_{\theta}$ )

outer expectation  
(against  $\mathbb{Q}$ )

# Nested expectations

Quantity of interest:

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta .$$

inner expectation  
(against  $\mathbb{P}_{\theta}$ )

outer expectation  
(against  $\mathbb{Q}$ )



When  $f$  is linear, it is a joint expectation, but we will usually be interested in **non-linear**  $f$ .

# Nested expectations

Quantity of interest:

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta .$$

inner expectation  
(against  $\mathbb{P}_{\theta}$ )

outer expectation  
(against  $\mathbb{Q}$ )



When  $f$  is linear, it is a joint expectation, but we will usually be interested in **non-linear**  $f$ .



I will be using only one level of nesting for simplicity, but we may sometimes care about more...

# Examples in Stats/ML/UQ

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$

# Examples in Stats/ML/UQ

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Batch active learning/Bayesian optimisation (inner: acquisition function for 1st point, outer: acquisition function for 2nd point given 1st point)



# Examples in Stats/ML/UQ

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Statistical divergences for conditional distributions (inner: standard statistical divergence, outer: average over conditioning variable)

# Examples in Stats/ML/UQ

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Bayesian distributionally robust optimisation (inner: expected risk against model, outer: expectation over posterior)

# Examples in Stats/ML/UQ

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Bayesian experimental design (inner: information gain - expectation over posterior, outer: expected information gain - expectation over marginal predictive distribution)

# Examples in other fields

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Option pricing (inner: expected loss given shock , outer: expectation over distribution of potential shocks)



# Examples in other fields

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



Health economics - Expected value of partial perfect information (inner: expected patient outcome given variable of interest, outer: expectation over prior beliefs about variable of interest)



# Examples in other fields

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta$$



# Some desiderata

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta.$$

- We define the **cost** of a method as the # function evaluations/samples needed for:

$$\left\{ \begin{array}{l} \text{Absolute error} = |I - \hat{I}| \leq \Delta \\ \text{RMSE} = \sqrt{\mathbb{E}[(I - \hat{I})^2]} \leq \Delta \end{array} \right.$$

# Some desiderata

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta.$$

- We define the **cost** of a method as the # function evaluations/samples needed for:

$$\left\{ \begin{array}{l} \text{Absolute error} = |I - \hat{I}| \leq \Delta \\ \text{RMSE} = \sqrt{\mathbb{E}[(I - \hat{I})^2]} \leq \Delta \end{array} \right.$$

- Ideally, we would like an estimator where

$$\text{Cost} = O(\Delta^{-r}) \text{ for (very) small } r$$



# Some desiderata

$$I := \int_{\Theta} f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right) q(\theta) d\theta.$$

- We define the **cost** of a method as the # function evaluations/samples needed for:

$$\left\{ \begin{array}{l} \text{Absolute error} = |I - \hat{I}| \leq \Delta \\ \text{RMSE} = \sqrt{\mathbb{E}[(I - \hat{I})^2]} \leq \Delta \end{array} \right.$$

- Ideally, we would like an estimator where

$$\text{Cost} = O(\Delta^{-r}) \text{ for (very) small } r$$



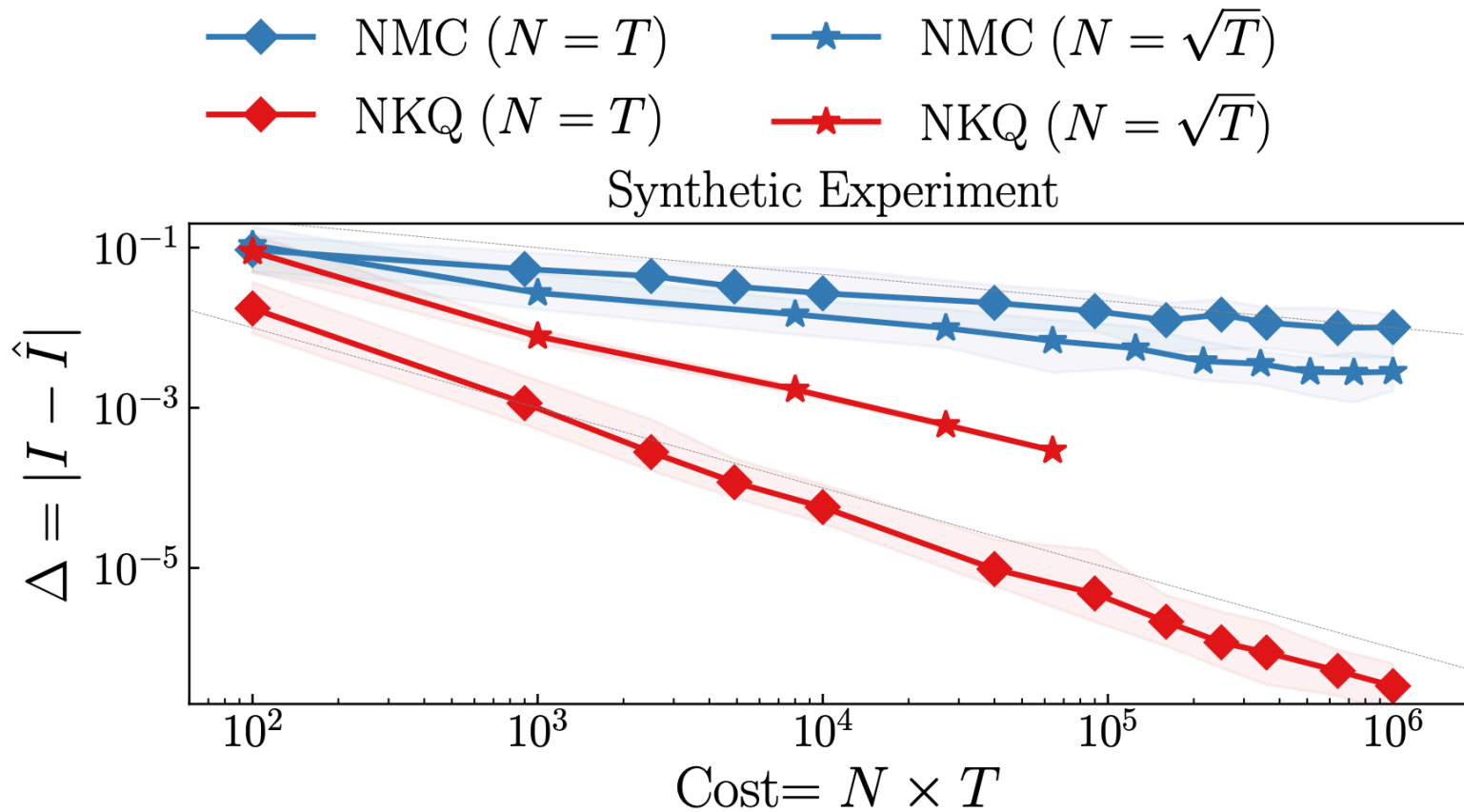
This is very important as most existing estimators tend to be **very** expensive.

# What to expect....

$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_{\theta} = U[0,1]$$



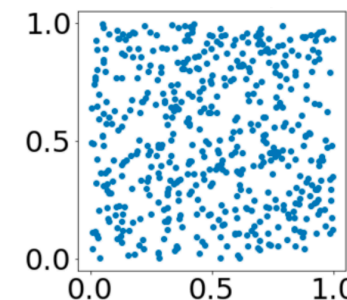
# Nested Monte Carlo

The most obvious estimator:

IID samples:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top \sim \mathbb{Q}$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$$



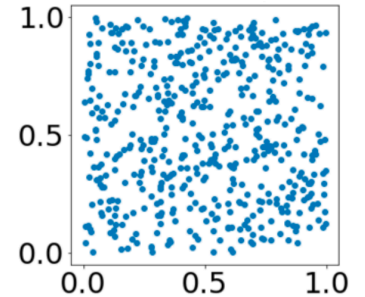
# Nested Monte Carlo

The most obvious estimator:

IID samples:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top \sim \mathbb{Q}$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$$



$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

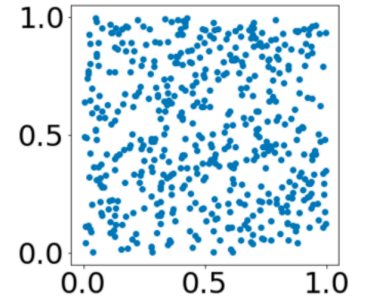
# Nested Monte Carlo

The most obvious estimator:

IID samples:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top \sim \mathbb{Q}$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$$



$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

inner Monte Carlo

outer Monte Carlo

# Convergence of Nested Monte Carlo

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  is Lipschitz, we get:  $\Delta \leq C_1 N^{-\frac{1}{2}} + C_2 T^{-\frac{1}{2}}$

# Convergence of Nested Monte Carlo

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  is Lipschitz, we get:

$$\Delta \leq C_1 N^{-\frac{1}{2}} + C_2 T^{-\frac{1}{2}}$$

Monte Carlo rate!



# Convergence of Nested Monte Carlo

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  is Lipschitz, we get:

$$\Delta \leq C_1 N^{-\frac{1}{2}} + C_2 T^{-\frac{1}{2}}$$

Monte Carlo rate!

- Taking  $N = T$  therefore leads to:

$$\text{Cost}(\hat{I}_{\text{NMC}}) = O(\Delta^{-4})$$



# Convergence of Nested Monte Carlo

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  is Lipschitz, we get:

$$\Delta \leq C_1 N^{-\frac{1}{2}} + C_2 T^{-\frac{1}{2}}$$

Monte Carlo rate!

- Taking  $N = T$  therefore leads to:

$$\text{Cost}(\hat{I}_{\text{NMC}}) = O(\Delta^{-4})$$

Too large...

# Convergence of Nested Monte Carlo

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  is Lipschitz, we get:

$$\Delta \leq C_1 N^{-\frac{1}{2}} + C_2 T^{-\frac{1}{2}}$$

Monte Carlo rate!

- Taking  $N = T$  therefore leads to:

$$\text{Cost}(\hat{I}_{\text{NMC}}) = O(\Delta^{-4})$$

Too large...

- Note:** This is biased since we never get to evaluate:

$$f \left( \int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx \right)$$

# Convergence of Nested Monte Carlo (ctd)

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

- Assuming  $f$  has bounded 2nd derivative , we get:

$$\Delta \leq C_1 N^{-1} + C_2 T^{-\frac{1}{2}}$$

# Convergence of Nested Monte Carlo (ctd)

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

Better than Monte Carlo rate!

- Assuming  $f$  has bounded 2nd derivative, we get:

$$\Delta \leq C_1 N^{-1} + C_2 T^{-\frac{1}{2}}$$

# Convergence of Nested Monte Carlo (ctd)

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

Better than Monte Carlo rate!

- Assuming  $f$  has bounded 2nd derivative, we get:

$$\Delta \leq C_1 N^{-1} + C_2 T^{-\frac{1}{2}}$$

- Taking  $N = \sqrt{T}$  leads to:

$$\text{Cost}(\hat{I}_{\text{NMC}}) = O(\Delta^{-3})$$

# Convergence of Nested Monte Carlo (ctd)

$$\hat{I}_{\text{NMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

Better than Monte Carlo rate!

- Assuming  $f$  has bounded 2nd derivative, we get:

$$\Delta \leq C_1 N^{-1} + C_2 T^{-\frac{1}{2}}$$

- Taking  $N = \sqrt{T}$  leads to:

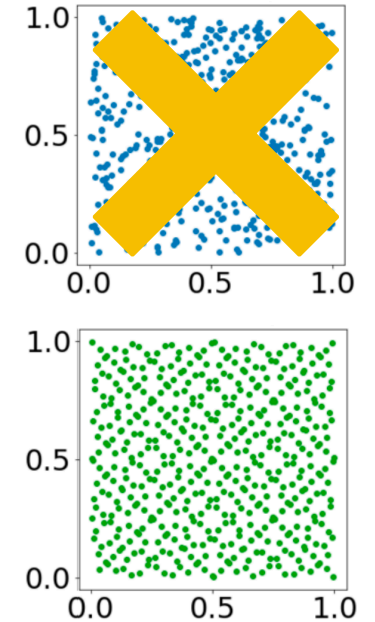
$$\text{Cost}(\hat{I}_{\text{NMC}}) = O(\Delta^{-3})$$

Smaller than 4, but still quite large

# Nested quasi Monte Carlo

QMC points:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$$
$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top, \quad t \in \{1, \dots, T\}$$



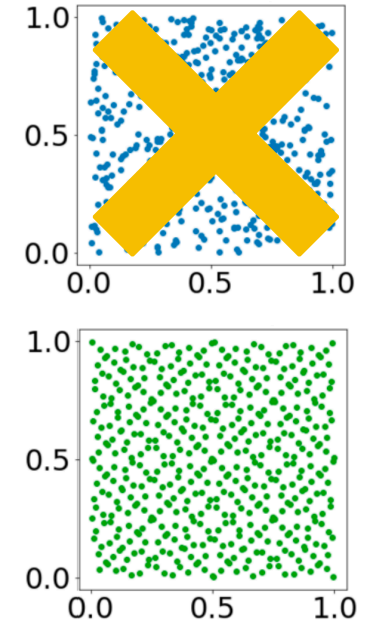
# Nested quasi Monte Carlo

QMC points:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top, \quad t \in \{1, \dots, T\}$$

$$\hat{I}_{\text{NQMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$





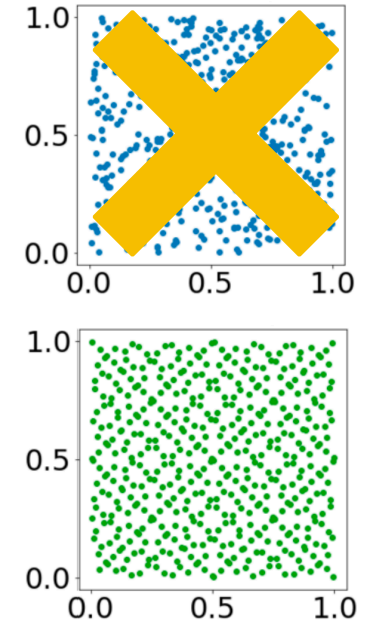
# Nested quasi Monte Carlo

QMC points:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top, \quad t \in \{1, \dots, T\}$$

$$\hat{I}_{\text{NQMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

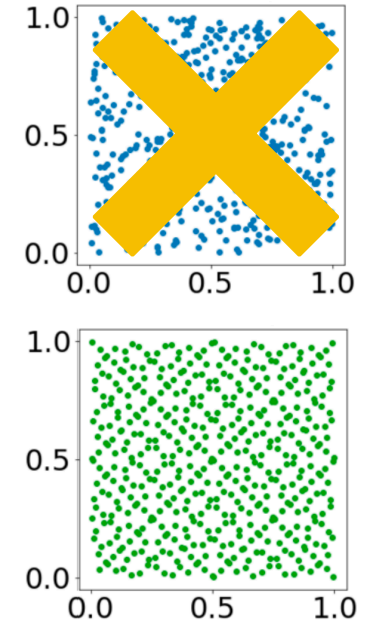


➔ Limited to cases where  $\mathcal{X} = [0,1]^{d_x}$ ,  $\Theta = [0,1]^{d_\theta}$  and uniform measures.

# Nested quasi Monte Carlo

QMC points:  $\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$   
 $x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top, \quad t \in \{1, \dots, T\}$

$$\hat{I}_{\text{NQMC}} := \frac{1}{T} \sum_{t=1}^T f \left( \frac{1}{N} \sum_{n=1}^N g(x_n^{(t)}, \theta_t) \right).$$

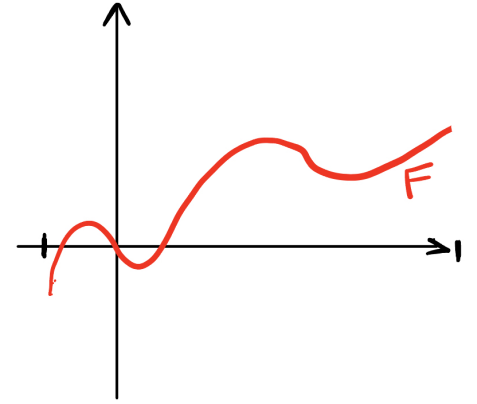


- ➔ Limited to cases where  $\mathcal{X} = [0,1]^{d_x}$ ,  $\Theta = [0,1]^{d_\theta}$  and uniform measures.
- ➔ Can get  $\text{Cost}(\hat{I}_{\text{NQMC}}) = O(\Delta^{-2.5})$  but requires **very strong** assumptions on  $f$  (second and third derivatives are monotone).

# Multi-level Monte Carlo

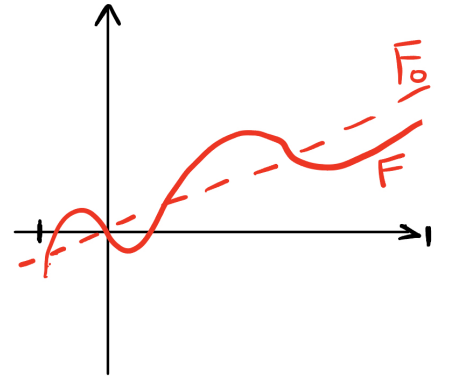
- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrand of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$I := \int_{\Theta} F(\theta) q(\theta) d\theta$$



# Multi-level Monte Carlo

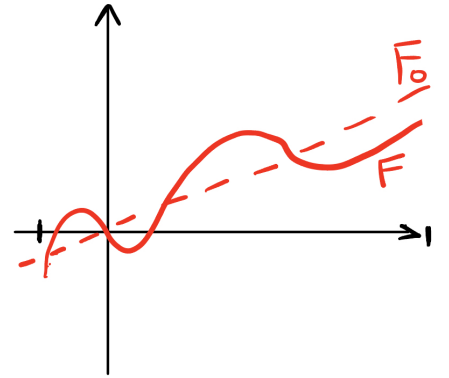
- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$



$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \int_{\Theta} (F(\theta) - F_0(\theta))q(\theta)d\theta$$

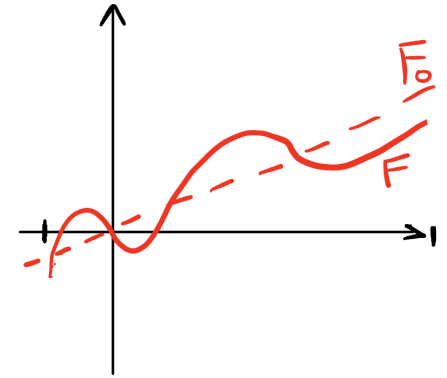
# Multi-level Monte Carlo

- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$



$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta$$

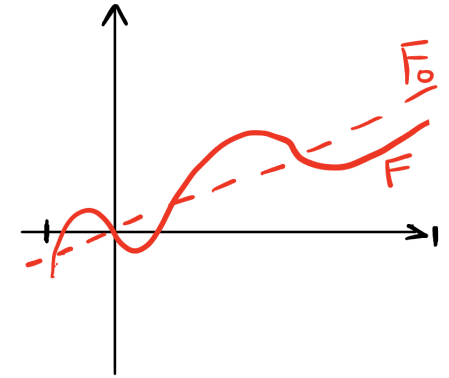
# Multi-level Monte Carlo



- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$\begin{aligned}
 I &:= \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta \\
 &\approx \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))
 \end{aligned}$$

# Multi-level Monte Carlo



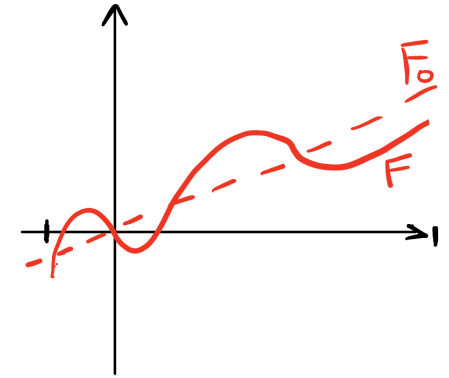
- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta$$

$$\approx \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

High variance

# Multi-level Monte Carlo



- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta$$

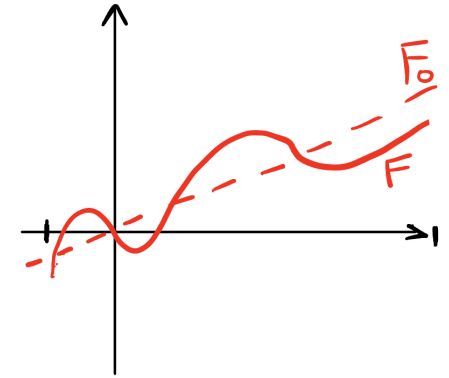
$$\approx \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

Large since cheap

High variance



# Multi-level Monte Carlo



- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta$$

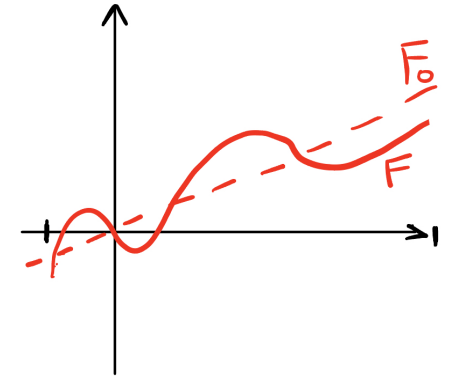
$$\approx \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

Large since cheap

High variance

(Very) small variance

# Multi-level Monte Carlo



- Back to IID points:  $\theta_{1:T}^l = (\theta_1, \dots, \theta_T)^\top, l \in \{0, \dots, L\}$
- Integrands of increasing fidelity and cost:  $F_0, F_1, \dots, F_L = F$

$$I := \int_{\Theta} F(\theta)q(\theta)d\theta = \int_{\Theta} F_0(\theta)q(\theta)d\theta + \sum_{l=1}^L \int_{\Theta} (F_l(\theta) - F_{l-1}(\theta))q(\theta)d\theta$$

$$\approx \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

Large since cheap

High variance

Small since expensive

(Very) small variance

# MLMC for nested expectations

IID points:  $\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$   
 $x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$

Define:  $F(\theta) := f\left(\int g(x, \theta) p_\theta(x) dx\right)$      $F_l(\theta) := f\left(\frac{1}{N_l} \sum_{n=1}^{N_l} g(x_n, \theta)\right)$

# MLMC for nested expectations

IID points:  $\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$   
 $x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$

Define:  $F(\theta) := f\left(\int g(x, \theta) p_\theta(x) dx\right) \quad F_l(\theta) := f\left(\frac{1}{N_l} \sum_{n=1}^{N_l} g(x_n, \theta)\right)$

$$\hat{I}_{\text{MLMC}} = \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

# MLMC for nested expectations

IID points:  $\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$   
 $x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$

Define:  $F(\theta) := f\left(\int g(x, \theta) p_\theta(x) dx\right) \quad F_l(\theta) := f\left(\frac{1}{N_l} \sum_{n=1}^{N_l} g(x_n, \theta)\right)$

$$\hat{I}_{\text{MLMC}} = \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$



Can get much lower cost:

$$\text{Cost}(\hat{I}_{\text{MLMC}}) = O(\Delta^{-2})$$

# MLMC for nested expectations

IID points:

$$\theta_{1:T} = (\theta_1, \dots, \theta_T)^\top$$

$$x_{1:N}^{(t)} = (x_1^{(t)}, \dots, x_N^{(t)})^\top \sim \mathbb{P}_{\theta_t}, \quad t \in \{1, \dots, T\}$$

Define:

$$F(\theta) := f \left( \int g(x, \theta) p_\theta(x) dx \right) \quad F_l(\theta) := f \left( \frac{1}{N_l} \sum_{n=1}^{N_l} g(x_n, \theta) \right)$$

$$\hat{I}_{\text{MLMC}} = \frac{1}{T_0} \sum_{t=1}^{T_0} F_0(\theta_t^0) + \sum_{l=1}^L \frac{1}{T_l} \sum_{t=1}^{T_l} (F_l(\theta_t^l) - F_{l-1}(\theta_t^l))$$

Much better than NMC/NQMC!



Can get much lower cost:

$$\text{Cost}(\hat{I}_{\text{MLMC}}) = O(\Delta^{-2})$$

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:



# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

$\mathcal{O}(10^5)$  evaluations

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

$\mathcal{O}(10^5)$  evaluations

$\mathcal{O}(10^4)$  evaluations

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$
<b>NKQ (Corollary 1)</b>	$\tilde{\mathcal{O}}\left(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}-\frac{d_{\Theta}}{s_{\Theta}}}\right)$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

$\mathcal{O}(10^5)$  evaluations

$\mathcal{O}(10^4)$  evaluations

( $\tilde{\mathcal{O}}$  means I am hiding log terms)

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$
<b>NKQ (Corollary 1)</b>	$\tilde{\mathcal{O}}\left(\Delta^{\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}} - \frac{d_{\Theta}}{s_{\Theta}}}\right)$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

$\mathcal{O}(10^5)$  evaluations

$\mathcal{O}(10^4)$  evaluations

( $\tilde{\mathcal{O}}$  means I am hiding log terms)

Smaller than 2 for sufficiently smooth integrands!

# A comparison of convergence rates

Method	Cost
NMC	$\mathcal{O}(\Delta^{-3})$ or $\mathcal{O}(\Delta^{-4})$
NQMC	$\mathcal{O}(\Delta^{-2.5})$
MLMC	$\mathcal{O}(\Delta^{-2})$
<b>NKQ (Corollary 1)</b>	$\tilde{\mathcal{O}}\left(\Delta^{\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}} - \frac{d_{\Theta}}{s_{\Theta}}}\right)$

To get  $\Delta = \mathcal{O}(0.01)$ , we need:

$\mathcal{O}(10^6)$  or  $\mathcal{O}(10^9)$  evaluations

$\mathcal{O}(10^5)$  evaluations

$\mathcal{O}(10^4)$  evaluations

←  $\mathcal{O}(10^2)$  or  $\mathcal{O}(10^3)$  evaluations?

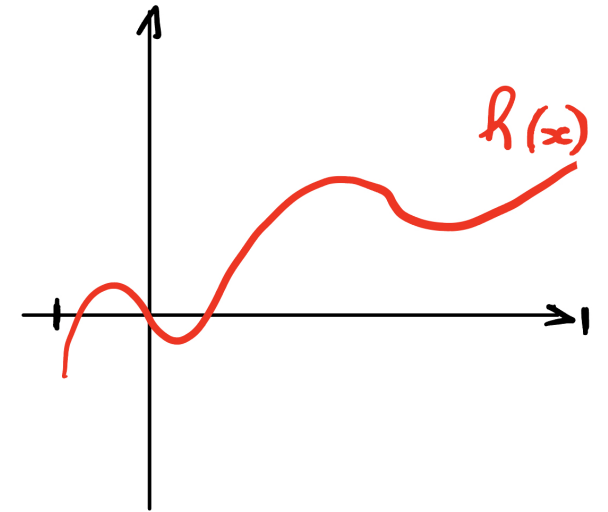
Smaller than 2 for sufficiently smooth integrands!

( $\tilde{\mathcal{O}}$  means I am hiding log terms)

# Quadrature rules

Quantity of interest:  $I = \int_{\mathcal{X}} h(x)\pi(x)dx$

Data:  $x_{1:N} := [x_1, \dots, x_N]^\top \in \mathcal{X}^N,$   
 $h(x_{1:N}) := [h(x_1), \dots, h(x_N)]^\top \in \mathbb{R}^N,$

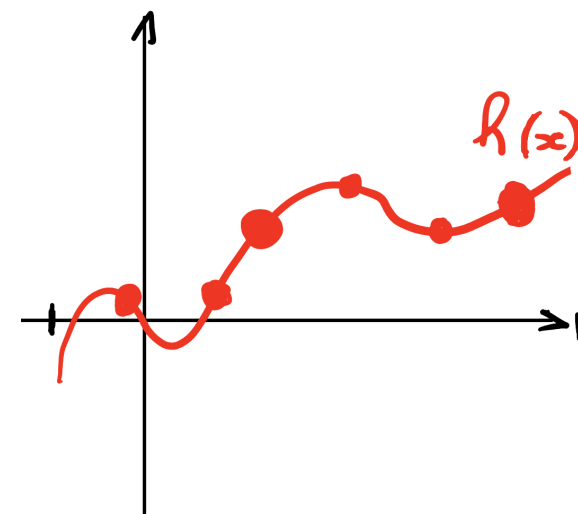


# Quadrature rules

Quantity of interest:  $I = \int_{\mathcal{X}} h(x)\pi(x)dx$

Data:  $x_{1:N} := [x_1, \dots, x_N]^\top \in \mathcal{X}^N,$   
 $h(x_{1:N}) := [h(x_1), \dots, h(x_N)]^\top \in \mathbb{R}^N,$

Quadrature rule:  $\hat{I} = \sum_{i=1}^N w_i h(x_i)$



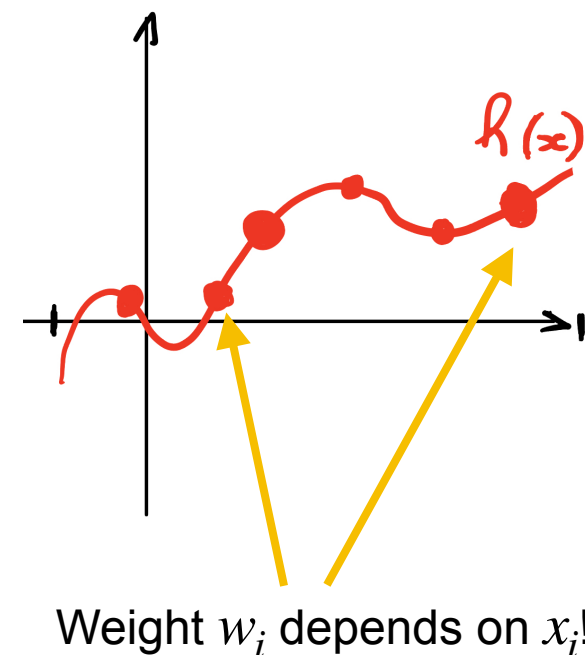


# Quadrature rules

Quantity of interest:  $I = \int_{\mathcal{X}} h(x)\pi(x)dx$

Data:  $x_{1:N} := [x_1, \dots, x_N]^\top \in \mathcal{X}^N,$   
 $h(x_{1:N}) := [h(x_1), \dots, h(x_N)]^\top \in \mathbb{R}^N,$

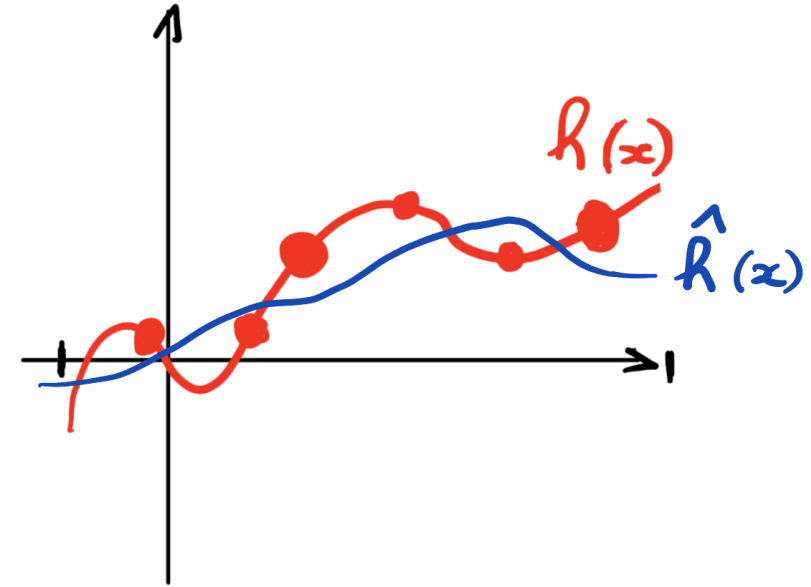
Quadrature rule:  $\hat{I} = \sum_{i=1}^N w_i h(x_i)$



# Kernel Quadrature (KQ)

- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

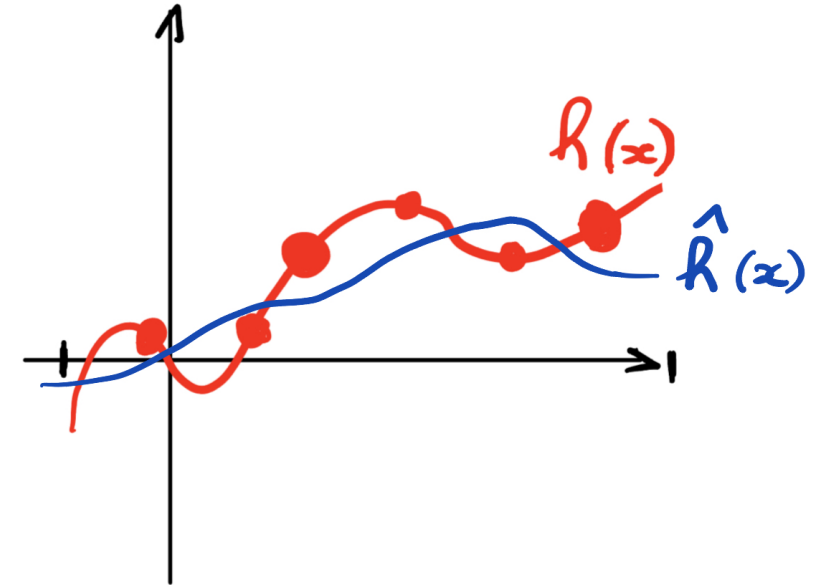


# Kernel Quadrature (KQ)

- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix



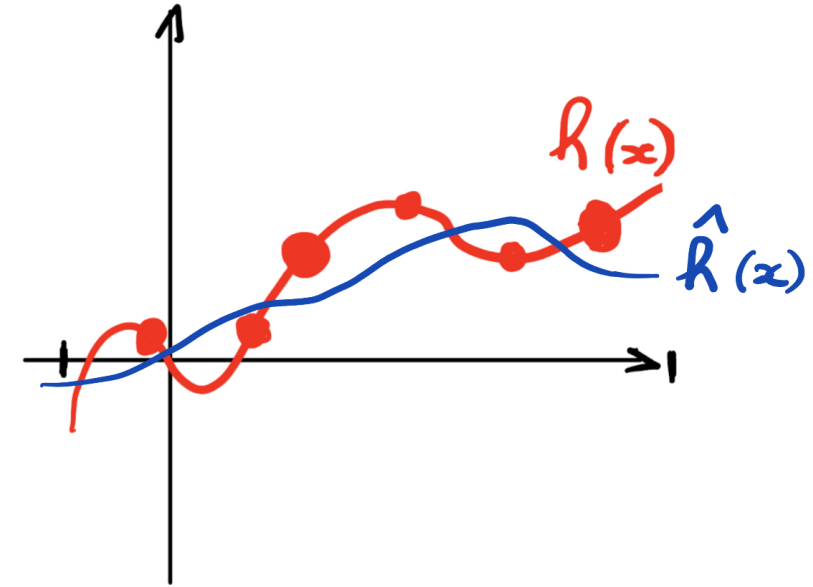
# Kernel Quadrature (KQ)

- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser



# Kernel Quadrature (KQ)

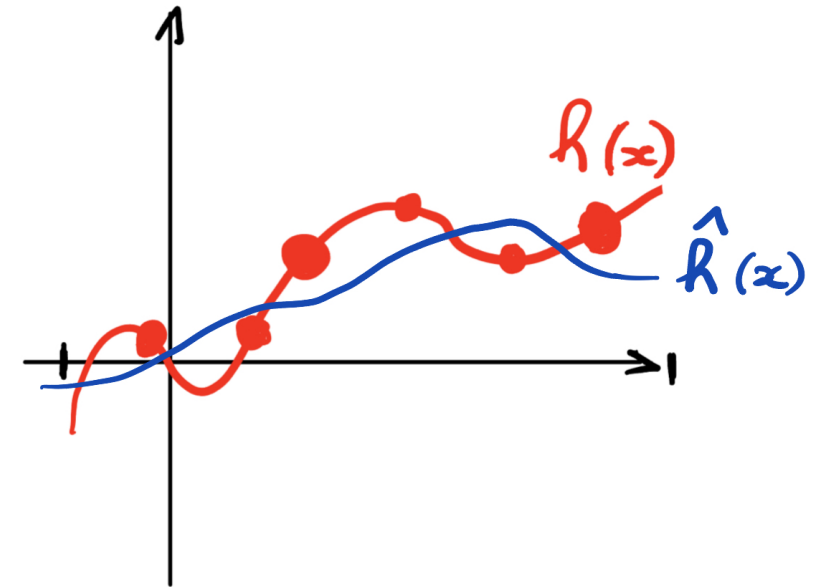
- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser

Identity matrix



# Kernel Quadrature (KQ)

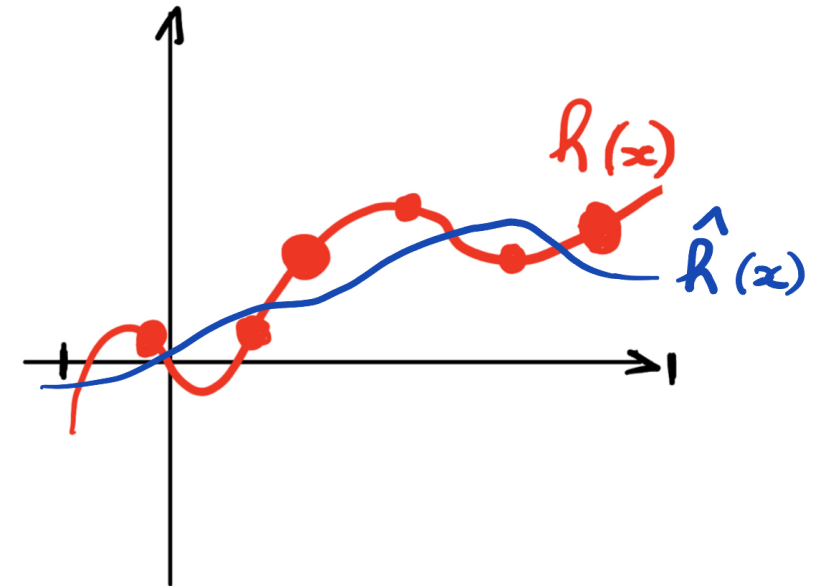
- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser

Identity matrix



- Use integral of  $\hat{h}$  as our estimator:

$$\hat{I}_{\text{KQ}} := \mu_{\pi}(x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N}) \quad \text{where } \mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

# Kernel Quadrature (KQ)

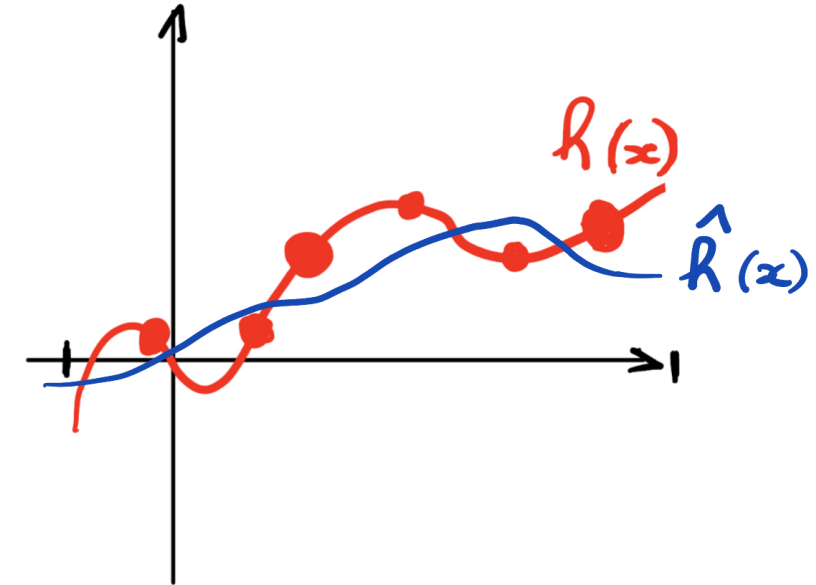
- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser

Identity matrix



- Use integral of  $\hat{h}$  as our estimator:

$$\hat{I}_{\text{KQ}} := \underbrace{\mu_{\pi}(x_{1:N})}_{\text{Weights!}} (k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1} h(x_{1:N}) \quad \text{where } \mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

# Kernel Quadrature (KQ)

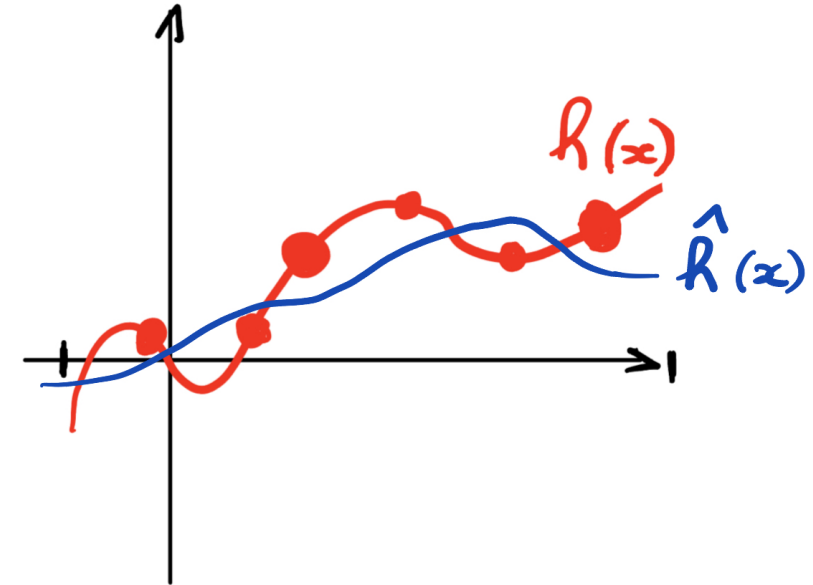
- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser

Identity matrix



- Use integral of  $\hat{h}$  as our estimator:

$$\hat{I}_{\text{KQ}} := \mu_{\pi}(x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Weights!

Kernel mean embedding!

$$\text{where } \mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$



# Kernel Quadrature (KQ)

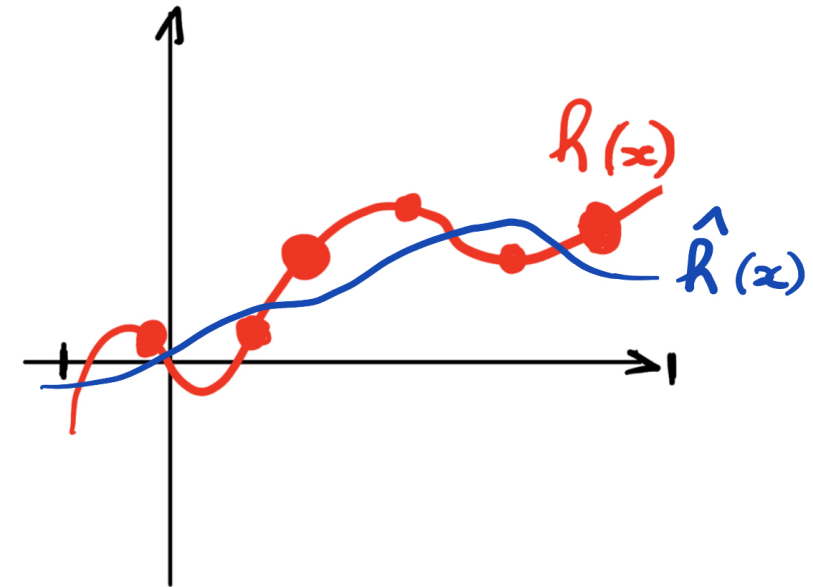
- Compute a kernel ridge regression estimator of  $h$  using some reproducing kernel  $k$ :

$$\hat{h}(x) := k(x, x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Gram matrix

Regulariser

Identity matrix



- Use integral of  $\hat{h}$  as our estimator:

$$\hat{I}_{\text{KQ}} := \mu_{\pi}(x_{1:N})(k(x_{1:N}, x_{1:N}) + N\lambda I_N)^{-1}h(x_{1:N})$$

Weights!

Kernel mean embedding!

where  $\mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$

- Closely relates to Bayesian quadrature (same procedure but with GP regression).

# Closed-form embeddings

We need closed-form kernel mean embeddings; this is not always straightforward!

$$\mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

# Closed-form embeddings

We need closed-form kernel mean embeddings; this is not always straightforward!

$$\mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

- **Trick 1: Stein reproducing kernels:** Construct a reproducing kernel only depending on  $\nabla_x \log \pi(x)$  where the kernel mean embedding is zero by construction!

# Closed-form embeddings

We need closed-form kernel mean embeddings; this is not always straightforward!

$$\mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

- **Trick 1: Stein reproducing kernels:** Construct a reproducing kernel only depending on  $\nabla_x \log \pi(x)$  where the kernel mean embedding is zero by construction!
- **Trick 2: change of variables:** Do a transformation so that the integral is against a simple measure for which we know a closed-form embedding formula.

# Closed-form embeddings

We need closed-form kernel mean embeddings; this is not always straightforward!

$$\mu_{\pi}(x) = \int_{\mathcal{X}} k(x, x')\pi(x')dx'$$

- **Trick 1: Stein reproducing kernels:** Construct a reproducing kernel only depending on  $\nabla_x \log \pi(x)$  where the kernel mean embedding is zero by construction!
- **Trick 2: change of variables:** Do a transformation so that the integral is against a simple measure for which we know a closed-form embedding formula.

# Advantages/Disadvantages of KQ

## Disadvantages:

- Computational cost is  $O(N^3)$  in the worst-case due to matrix inversion.
- Need closed-form kernel mean embeddings (but can be mitigated with two tricks).

# Advantages/Disadvantages of KQ

## Disadvantages:

- Computational cost is  $O(N^3)$  in the worst-case due to matrix inversion.
- Need closed-form kernel mean embeddings (but can be mitigated with two tricks).

## Advantages:

- Typically converges much faster than alternative estimators when integrand is smooth and not too high-dimensional!

# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

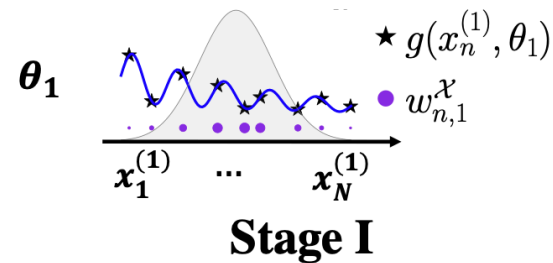
Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .



# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

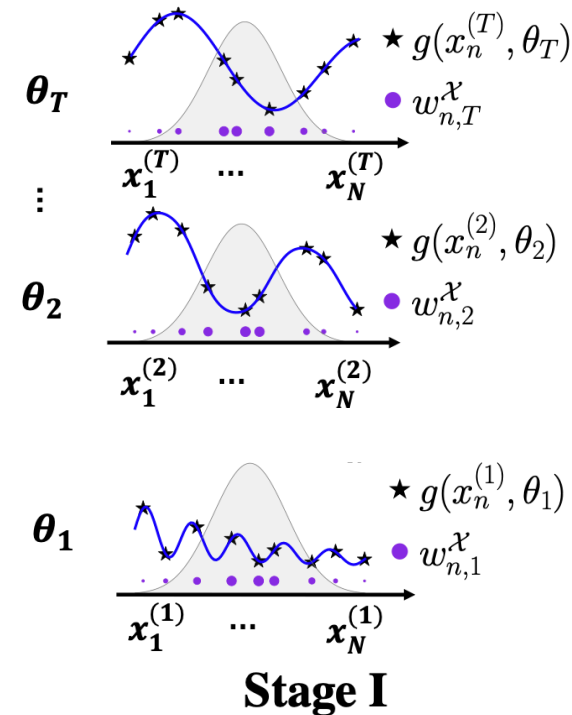
Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .



# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .



# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

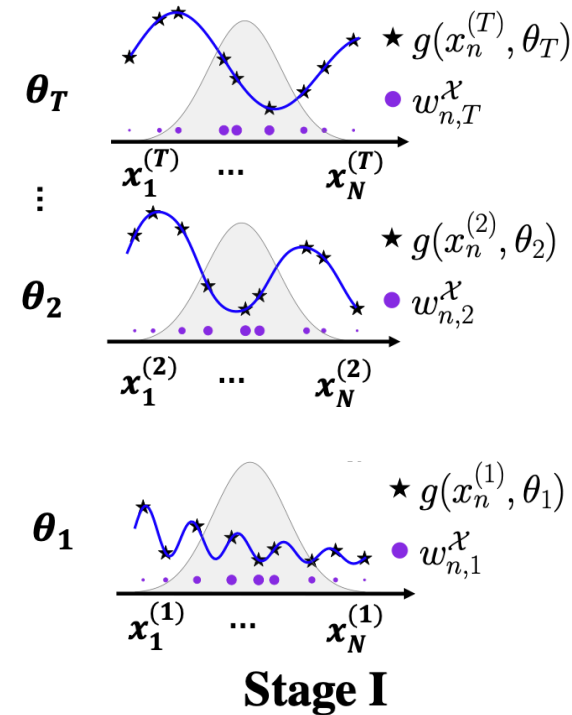
Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .

**Stage II:** Compute a KQ estimator (with  $k_{\Theta}$ ) for outer expectation; i.e. integral of

$$F(\theta) = f\left(\int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx\right)$$

using noisy data from stage I:

$$\hat{F}_{\text{KQ}}(\theta_t) = f(\hat{J}_{\text{KQ}}(\theta_t)) \approx F(\theta_t)$$



# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

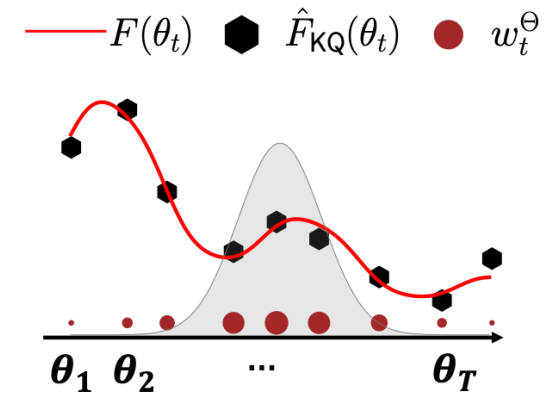
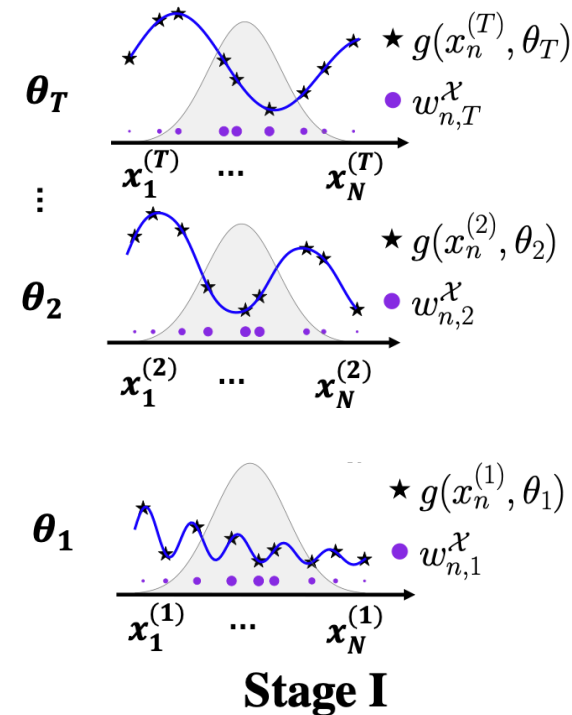
Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .

**Stage II:** Compute a KQ estimator (with  $k_{\Theta}$ ) for outer expectation; i.e. integral of

$$F(\theta) = f\left(\int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx\right)$$

using noisy data from stage I:

$$\hat{F}_{\text{KQ}}(\theta_t) = f(\hat{J}_{\text{KQ}}(\theta_t)) \approx F(\theta_t)$$



Stage II

# Nested Kernel Quadrature

**Stage I:** Compute KQ estimators (with  $k_{\mathcal{X}}$ ) for inner expectations; i.e. for integrals of  $g(\cdot, \theta_1), \dots, g(\cdot, \theta_T)$ .

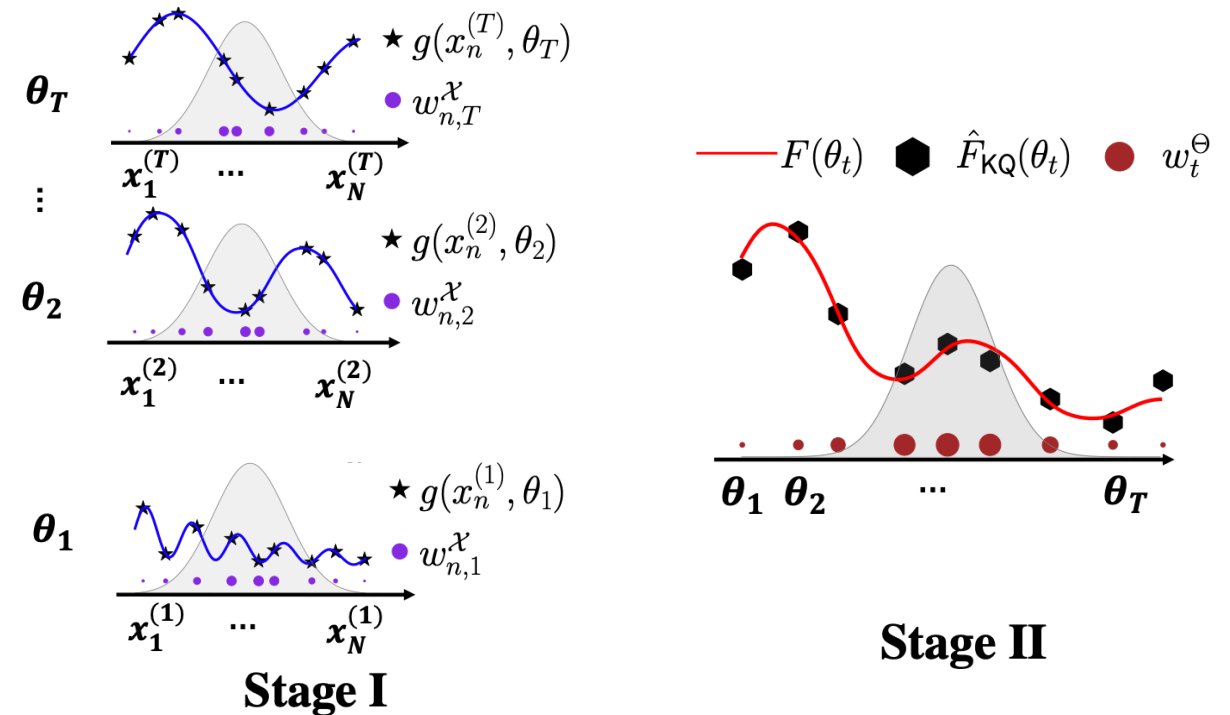
Denote these  $\hat{J}_{\text{KQ}}(\theta_1), \dots, \hat{J}_{\text{KQ}}(\theta_T)$ .

**Stage II:** Compute a KQ estimator (with  $k_{\Theta}$ ) for outer expectation; i.e. integral of

$$F(\theta) = f\left(\int_{\mathcal{X}} g(x, \theta) p_{\theta}(x) dx\right)$$

using noisy data from stage I:

$$\hat{F}_{\text{KQ}}(\theta_t) = f(\hat{J}_{\text{KQ}}(\theta_t)) \approx F(\theta_t)$$



$$\hat{I}_{\text{NKQ}} = \sum_{t=1}^T w_t^{\Theta} f\left(\sum_{i=1}^N w_{n,t}^{\mathcal{X}} g(x_n^{(t)}, \theta_t)\right)$$

# Convergence guarantees for NKQ

- **Theorem (informal):** Let  $\mathcal{X} = [0,1]^{d_x}$  and  $\Theta = [0,1]^{d_\Theta}$ . Under regularity assumptions including

# Convergence guarantees for NKQ

- **Theorem (informal):** Let  $\mathcal{X} = [0,1]^{d_x}$  and  $\Theta = [0,1]^{d_\theta}$ . Under regularity assumptions including
  - The samples  $x_{1:N}^{(t)}$  and  $\theta_{1:T}$  are iid from  $\mathbb{P}_{\theta_t}$  and  $\mathbb{Q}$  respectively.

# Convergence guarantees for NKQ

- **Theorem (informal):** Let  $\mathcal{X} = [0,1]^{d_{\mathcal{X}}}$  and  $\Theta = [0,1]^{d_{\Theta}}$ . Under regularity assumptions including
  - The samples  $x_{1:N}^{(t)}$  and  $\theta_{1:T}$  are iid from  $\mathbb{P}_{\theta_t}$  and  $\mathbb{Q}$  respectively.
  - The kernels  $k_{\mathcal{X}}$  and  $k_{\Theta}$  have (Sobolev) smoothness  $s_{\mathcal{X}} \geq \frac{d_{\mathcal{X}}}{2}$  and  $s_{\Theta} \geq \frac{d_{\Theta}}{2}$  respectively.



# Convergence guarantees for NKQ

- **Theorem (informal):** Let  $\mathcal{X} = [0,1]^{d_{\mathcal{X}}}$  and  $\Theta = [0,1]^{d_{\Theta}}$ . Under regularity assumptions including
  - The samples  $x_{1:N}^{(t)}$  and  $\theta_{1:T}$  are iid from  $\mathbb{P}_{\theta_t}$  and  $\mathbb{Q}$  respectively.
  - The kernels  $k_{\mathcal{X}}$  and  $k_{\Theta}$  have (Sobolev) smoothness  $s_{\mathcal{X}} \geq \frac{d_{\mathcal{X}}}{2}$  and  $s_{\Theta} \geq \frac{d_{\Theta}}{2}$  respectively.
  - $f \in C_b^{s_{\Theta}+1}$ ,  $\theta \mapsto p_{\theta}(x)$  and  $\theta \mapsto g(x, \theta)$  have smoothness  $s_{\Theta}$ .

# Convergence guarantees for NKQ

- **Theorem (informal):** Let  $\mathcal{X} = [0,1]^{d_x}$  and  $\Theta = [0,1]^{d_\Theta}$ . Under regularity assumptions including
  - The samples  $x_{1:N}^{(t)}$  and  $\theta_{1:T}$  are iid from  $\mathbb{P}_{\theta_t}$  and  $\mathbb{Q}$  respectively.
  - The kernels  $k_{\mathcal{X}}$  and  $k_{\Theta}$  have (Sobolev) smoothness  $s_{\mathcal{X}} \geq \frac{d_x}{2}$  and  $s_{\Theta} \geq \frac{d_{\Theta}}{2}$  respectively.
  - $f \in C_b^{s_{\Theta}+1}$ ,  $\theta \mapsto p_{\theta}(x)$  and  $\theta \mapsto g(x, \theta)$  have smoothness  $s_{\Theta}$ .
  - $x \mapsto D_{\theta}^{\beta} g(x, \theta)$  has smoothness  $s_{\mathcal{X}}$  for all  $\beta \in \mathbb{N}_0^{d_{\Theta}}$  with  $|\beta| \leq s_{\Theta}$ .

# Convergence guarantees for NKQ

Then, taking  $\lambda_{\mathcal{X}} = \tilde{O}(N^{-\frac{2s_{\mathcal{X}}}{d_{\mathcal{X}}}})$  and  $\lambda_{\Theta} = \tilde{O}(T^{-\frac{2s_{\Theta}}{d_{\Theta}}})$ , we get that for  $N, T$  large enough the following holds with high prob:

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O}\left(N^{-\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{-\frac{s_{\Theta}}{d_{\Theta}}}\right)$$

# Convergence guarantees for NKQ

Then, taking  $\lambda_{\mathcal{X}} = \tilde{O}(N^{-\frac{2s_{\mathcal{X}}}{d_{\mathcal{X}}}})$  and  $\lambda_{\Theta} = \tilde{O}(T^{-\frac{2s_{\Theta}}{d_{\Theta}}})$ , we get that for  $N, T$  large enough the following holds with high prob:

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O}\left( N^{\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{\frac{s_{\Theta}}{d_{\Theta}}} \right)$$

Fast! Much better than MC!

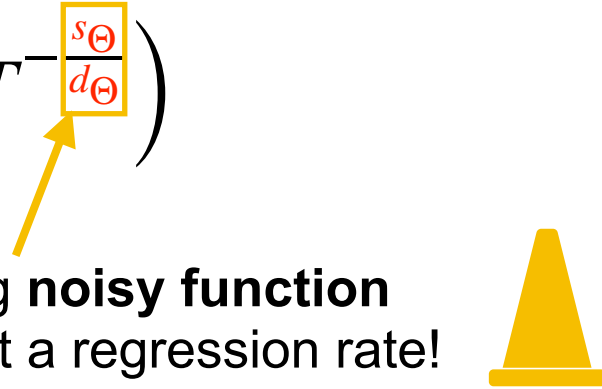
# Cost of NKQ

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O} \left( N^{-\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{-\frac{s_{\Theta}}{d_{\Theta}}} \right)$$

- This **fast (interpolation-type) rate is surprising** given we are using **noisy function values** for approximating the outer expectation, so we should expect a regression rate!



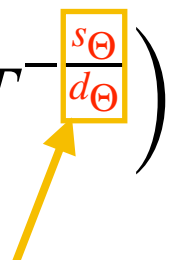
# Cost of NKQ

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O} \left( N^{-\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{-\frac{s_{\Theta}}{d_{\Theta}}} \right)$$


- This **fast (interpolation-type) rate is surprising** given we are using **noisy function values** for approximating the outer expectation, so we should expect a regression rate!

- Taking  $N = \tilde{O}(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}})$  and  $T = \tilde{O}(\Delta^{-\frac{d_{\Theta}}{s_{\Theta}}})$ , we get:  $\text{Cost}(\hat{I}_{\text{NKQ}}) = \tilde{O}(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}} - \frac{d_{\Theta}}{s_{\Theta}}})$

# Cost of NKQ

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O} \left( N^{-\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{-\frac{s_{\Theta}}{d_{\Theta}}} \right)$$


- This **fast (interpolation-type) rate is surprising** given we are using **noisy function values** for approximating the outer expectation, so we should expect a regression rate! 

- Taking  $N = \tilde{O}(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}})$  and  $T = \tilde{O}(\Delta^{-\frac{d_{\Theta}}{s_{\Theta}}})$ , we get:  $\text{Cost}(\hat{I}_{\text{NKQ}}) = \tilde{O}(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}} - \frac{d_{\Theta}}{s_{\Theta}}})$

- Recall that  $s_{\mathcal{X}} \geq d_{\mathcal{X}}/2$  and  $s_{\Theta} \geq d_{\Theta}/2$ , so we get at worst the NMC rate:  $\tilde{O}(\Delta^{-4})$

# Cost of NKQ

$$\Delta = \left| \hat{I}_{\text{NKQ}} - I \right| = \tilde{O}\left(N^{-\frac{s_{\mathcal{X}}}{d_{\mathcal{X}}}} + T^{-\frac{s_{\Theta}}{d_{\Theta}}}\right)$$

→ When  $s_{\mathcal{X}}$  and  $s_{\Theta}$  are large enough, we can beat NQMC and MLMC!



# Back to the synthetic experiment

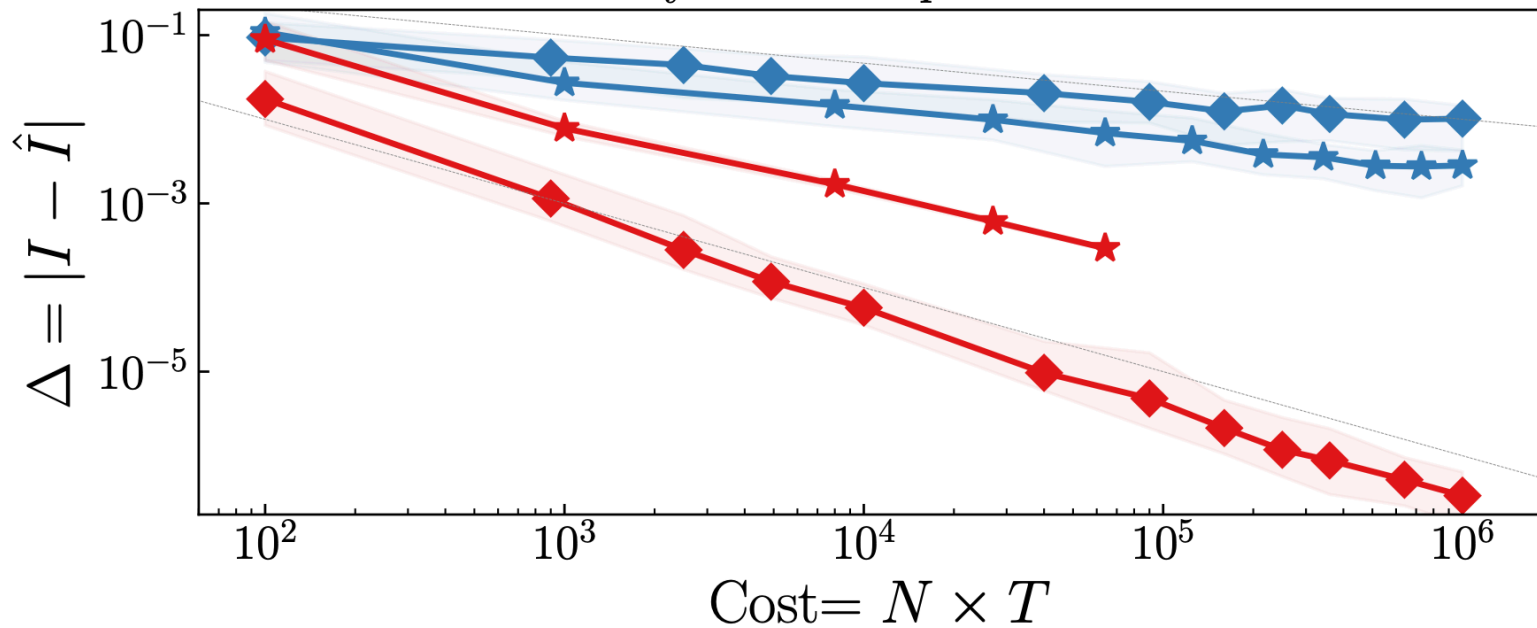
$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

- ◆— NMC ( $N = T$ )      —★— NMC ( $N = \sqrt{T}$ )
- ◆— NKQ ( $N = T$ )      —★— NKQ ( $N = \sqrt{T}$ )

Synthetic Experiment



# Back to the synthetic experiment

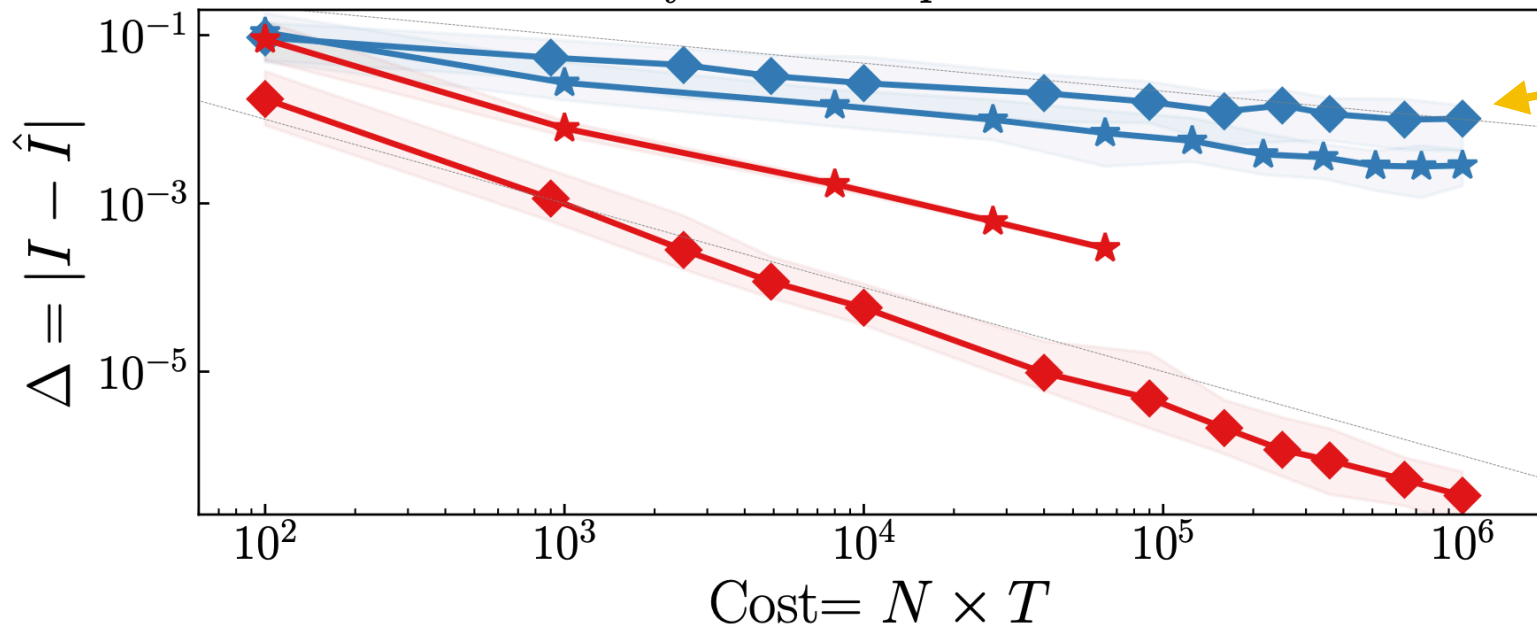
$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

- ◆— NMC ( $N = T$ )      —★— NMC ( $N = \sqrt{T}$ )
- ◆— NKQ ( $N = T$ )      —★— NKQ ( $N = \sqrt{T}$ )

Synthetic Experiment



As predicted by NMC theory,  $N = T$  gives the slower rate:  $O(\Delta^{-4})$

# Back to the synthetic experiment

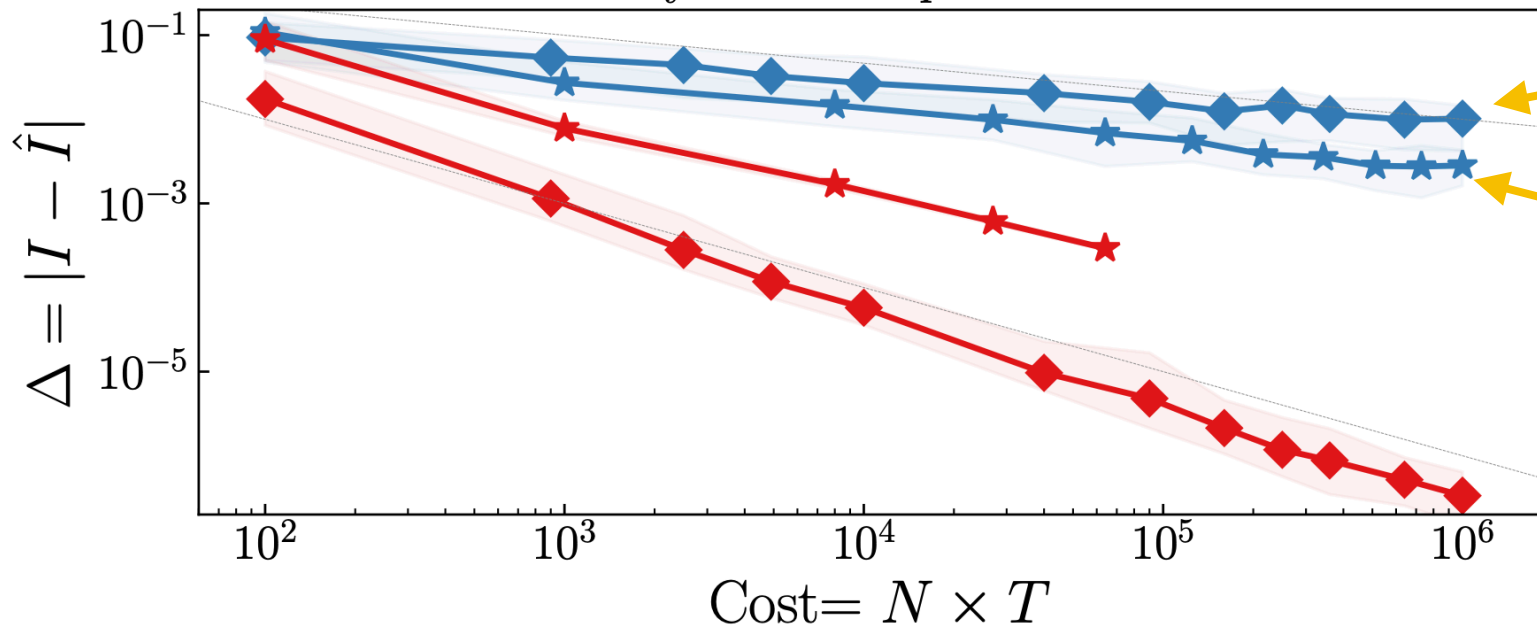
$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

- ◆— NMC ( $N = T$ )      —★— NMC ( $N = \sqrt{T}$ )
- ◆— NKQ ( $N = T$ )      —★— NKQ ( $N = \sqrt{T}$ )

Synthetic Experiment



As predicted by NMC theory,  $N = T$  gives the slower rate:  $O(\Delta^{-4})$

Taking  $N = \sqrt{T}$  gives the faster rate:  $O(\Delta^{-3})$  since integrand nice.

# Back to the synthetic experiment

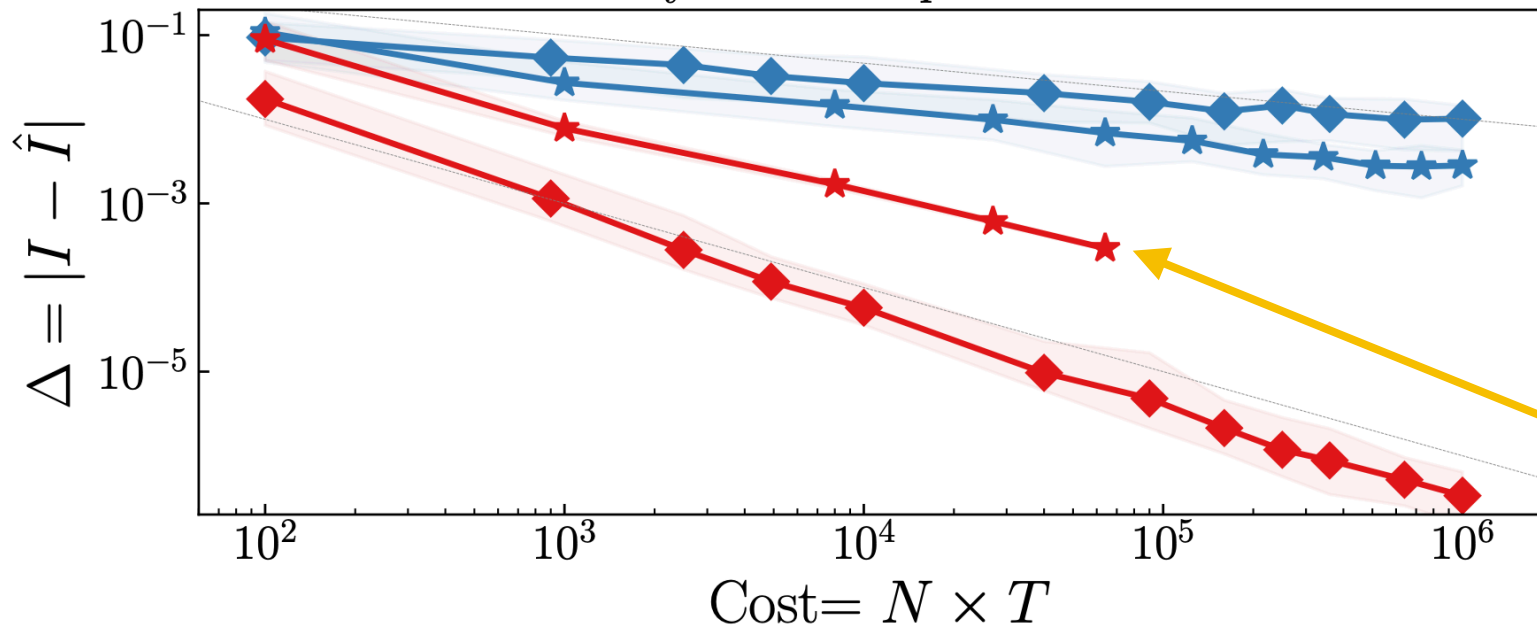
$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

- ◆— NMC ( $N = T$ )      —★— NMC ( $N = \sqrt{T}$ )
- ◆— NKQ ( $N = T$ )      —★— NKQ ( $N = \sqrt{T}$ )

Synthetic Experiment



For NKQ,  $N = \sqrt{T}$  is sub-optimal.

# Back to the synthetic experiment

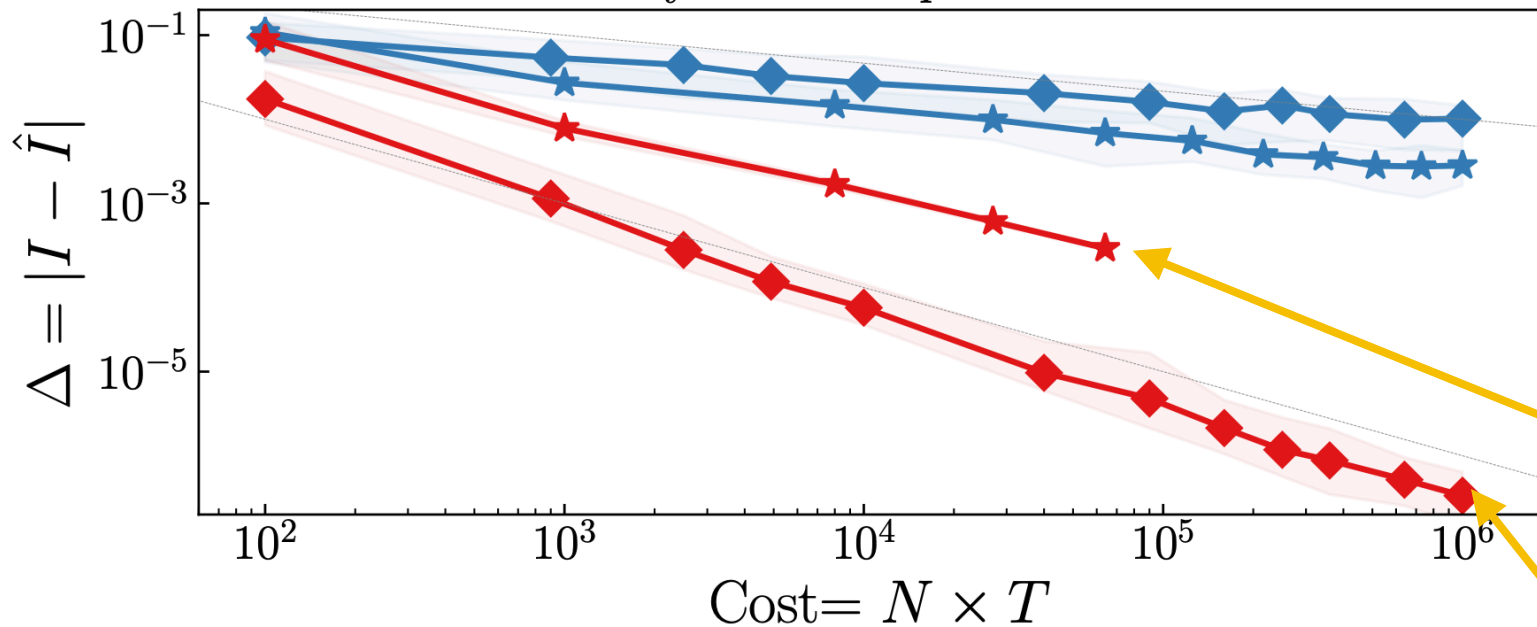
$$g(x, \theta) = x^{\frac{5}{2}} + \theta^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

- ◆— NMC ( $N = T$ )      —★— NMC ( $N = \sqrt{T}$ )
- ◆— NKQ ( $N = T$ )      —★— NKQ ( $N = \sqrt{T}$ )

Synthetic Experiment



For NKQ,  $N = \sqrt{T}$  is sub-optimal.

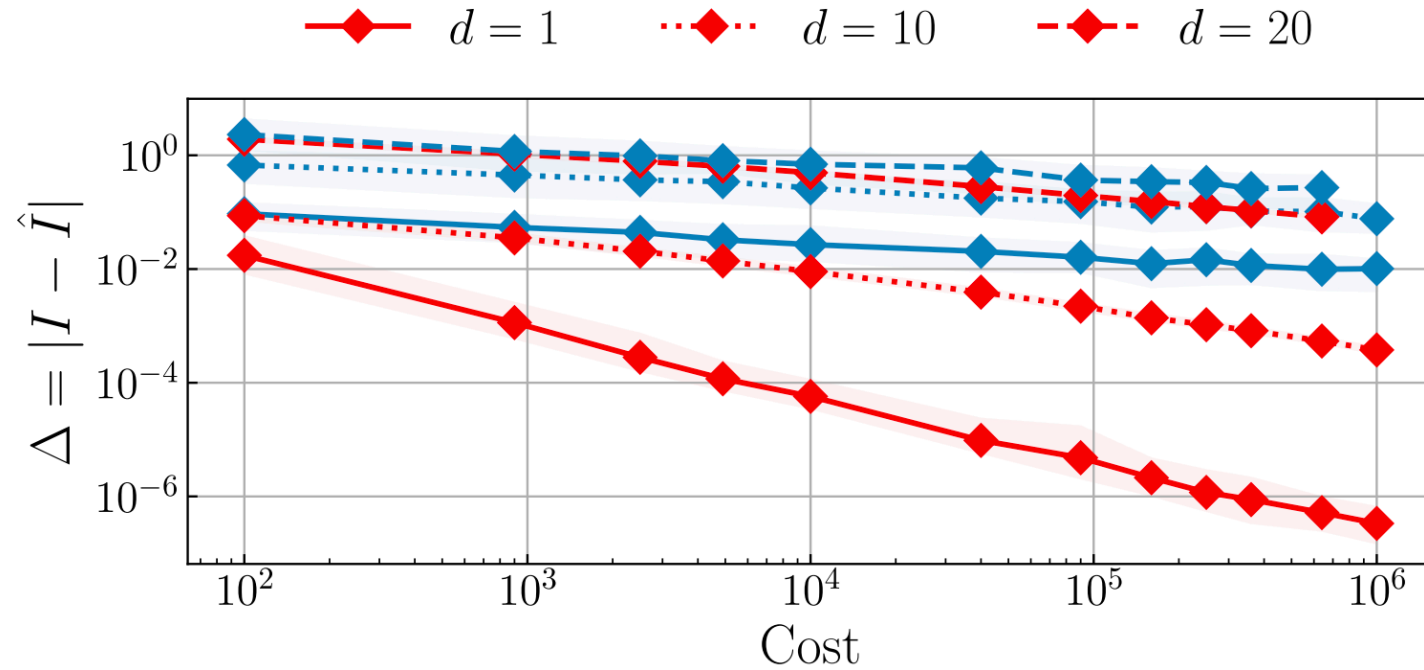
As predicted by our theorem,  $N = T$  gives a fast rate!

# We suffer in high dimensions...

$$g(x, \theta) = \|x\|_2^{\frac{5}{2}} + \|\theta\|_2^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

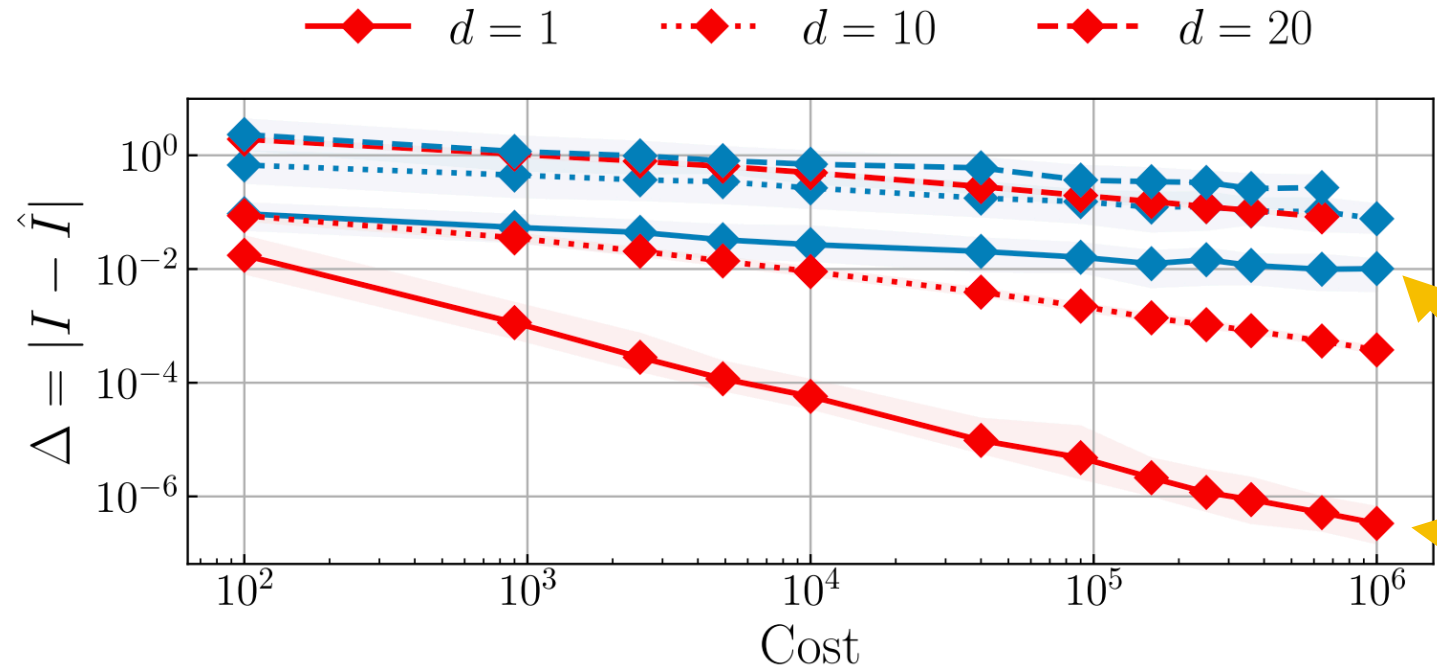


# We suffer in high dimensions...

$$g(x, \theta) = \|x\|_2^{\frac{5}{2}} + \|\theta\|_2^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$



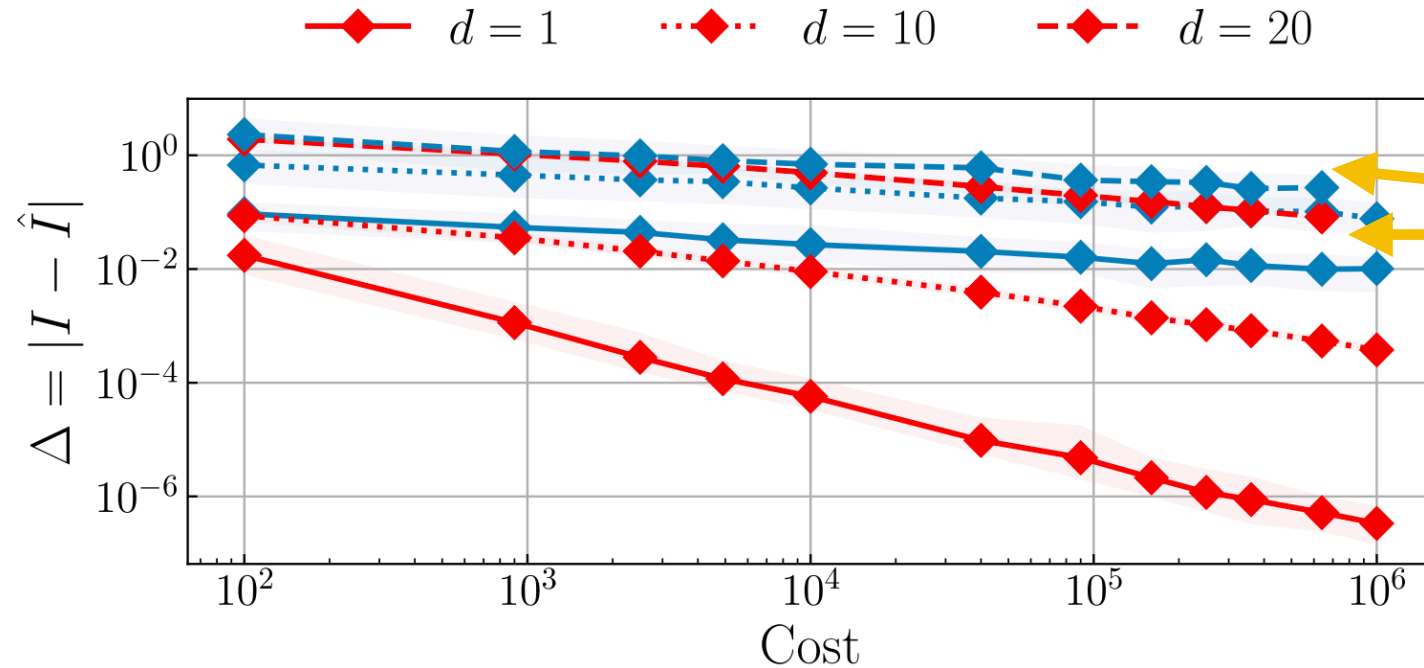
We do great in one dimension (as seen previously)

# We suffer in high dimensions...

$$g(x, \theta) = \|x\|_2^{\frac{5}{2}} + \|\theta\|_2^{\frac{5}{2}}$$

$$f(z) = z^2$$

$$\mathbb{Q} = \mathbb{P}_\theta = U[0,1]$$

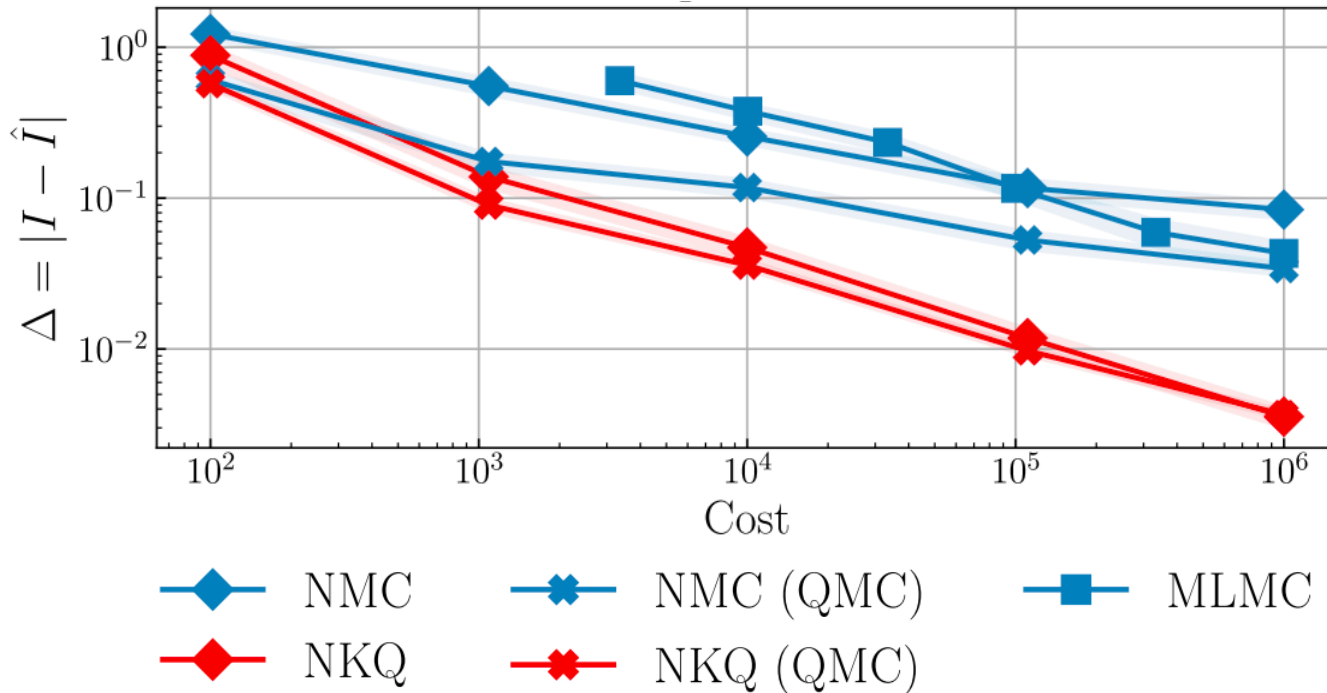


We do much worse in higher dimensions!

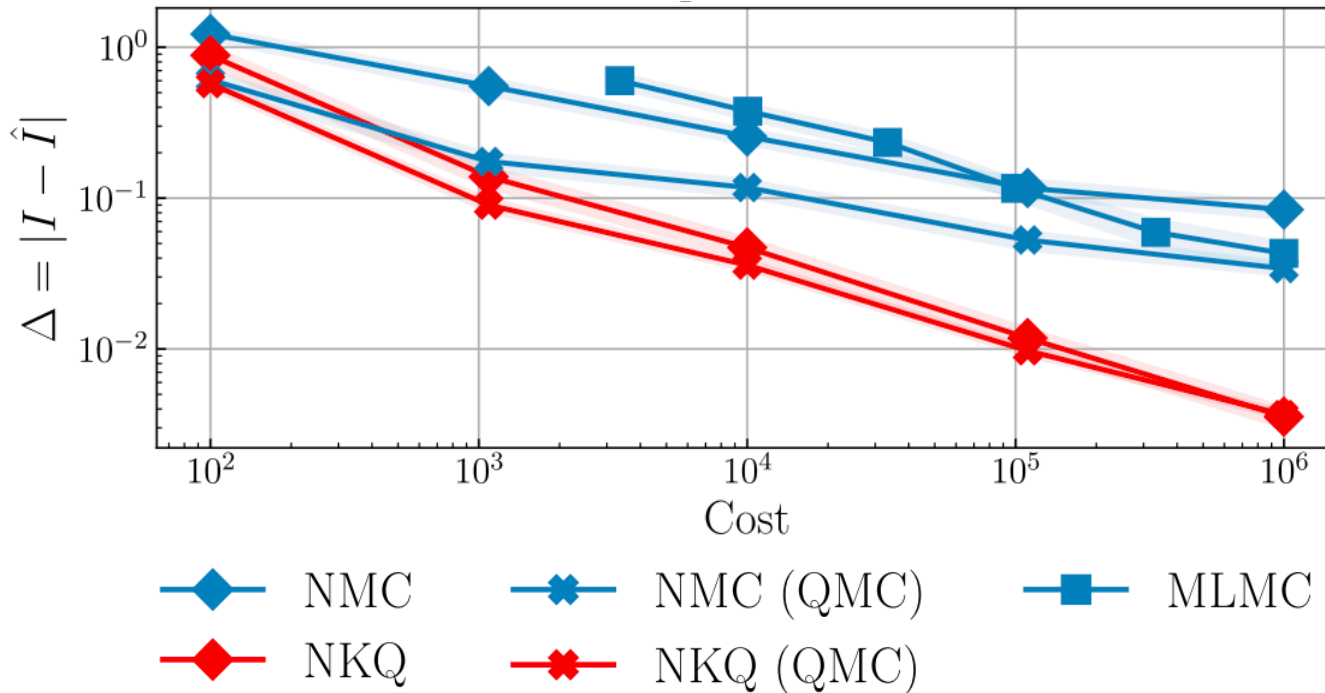


# Option pricing

- **Problem:** Pricing of options; expected loss of portfolio in the presence of potential economic shock.



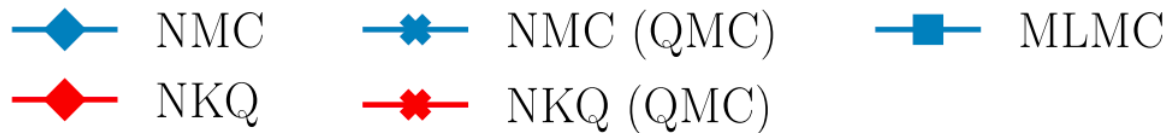
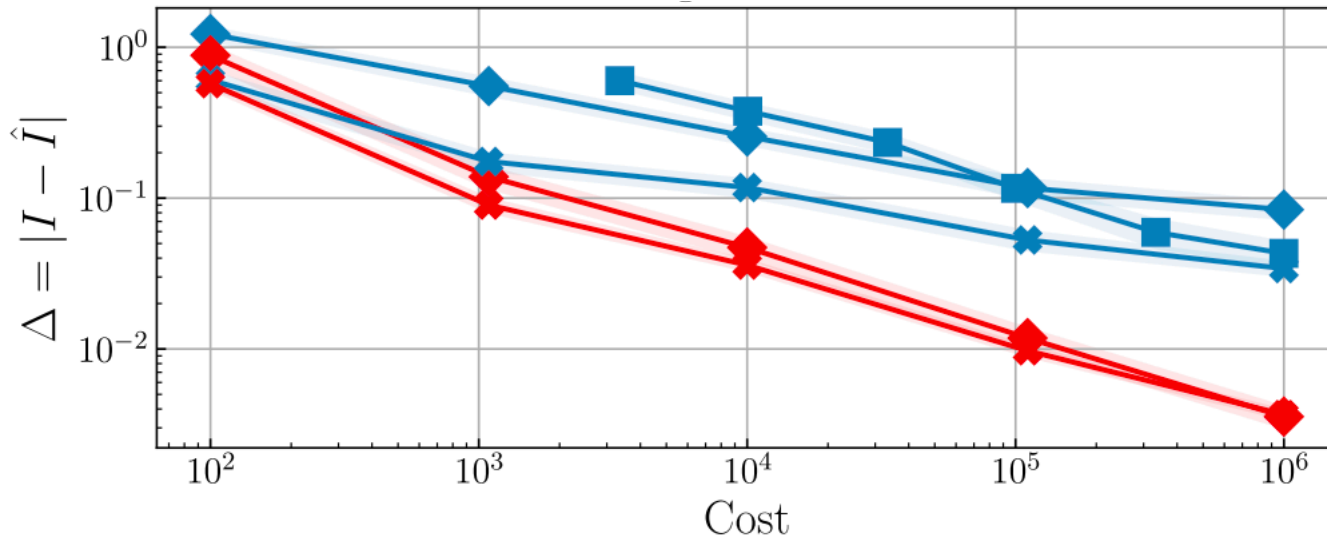
# Option pricing



- **Problem:** Pricing of options; expected loss of portfolio in the presence of potential economic shock.
- Some of the assumptions are broken (unbounded domain,  $f \notin C_b^2$ )
- We used  $k_{\mathcal{X}}$  and  $k_{\Theta}$  with  $s_{\mathcal{X}} = s_{\Theta} = 1$ .

# Option pricing

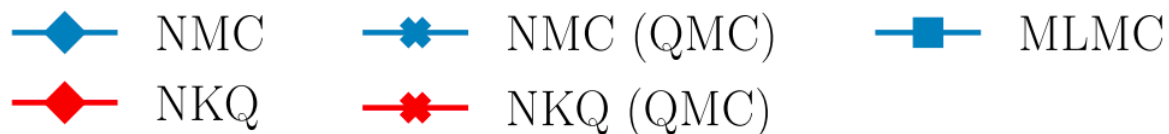
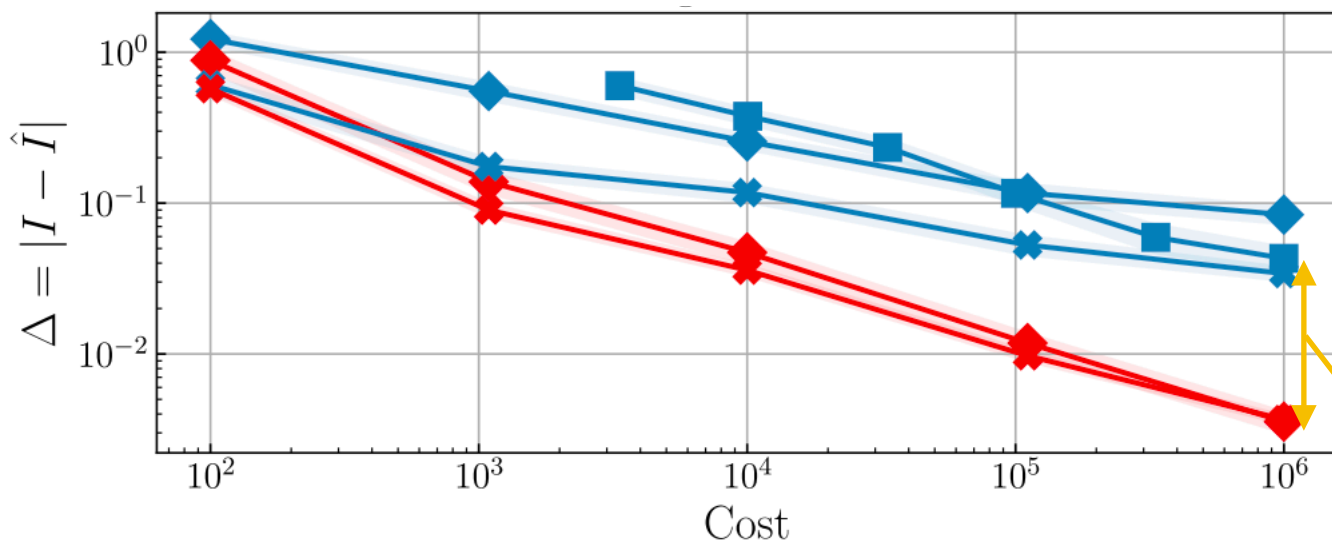
- **Problem:** Pricing of options; expected loss of portfolio in the presence of potential economic shock.



Approx **100x smaller** error than NMC

# Option pricing

- **Problem:** Pricing of options; expected loss of portfolio in the presence of potential economic shock.

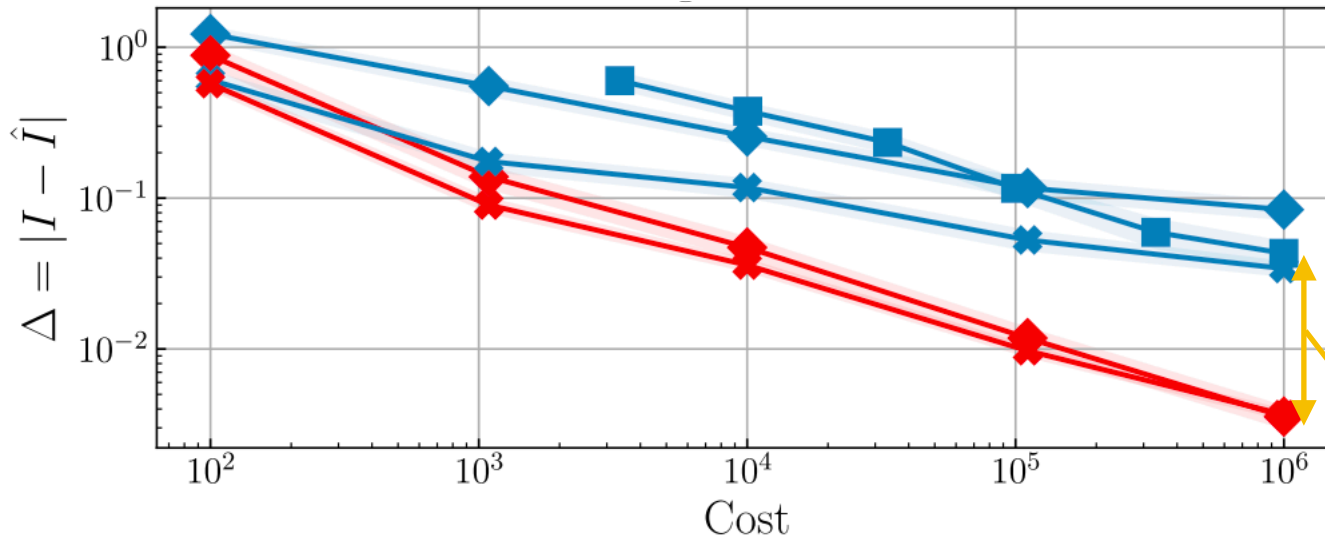


Approx **100x smaller** error than NMC

Approx **10x smaller** error than MLMC

# Option pricing

- **Problem:** Pricing of options; expected loss of portfolio in the presence of potential economic shock.



- ◆ NMC
- ◆ NMC (QMC)
- MLMC
- ◆ NKQ
- ◆ NKQ (QMC)

Approx **100x smaller** error than NMC

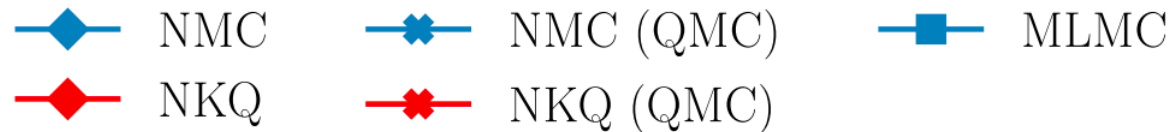
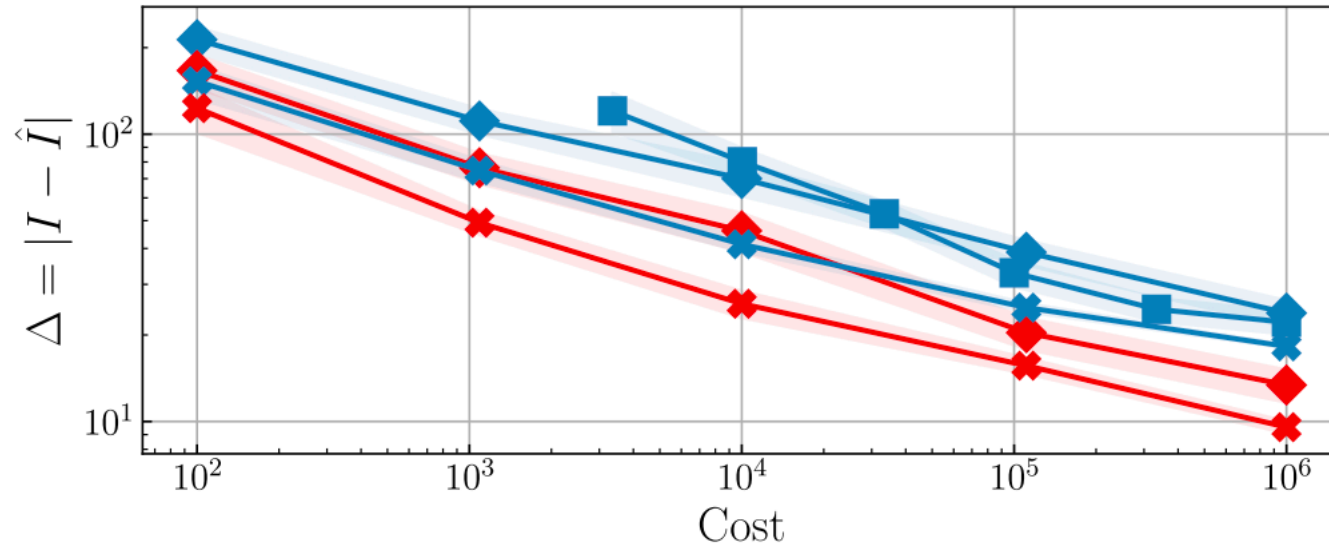
Approx **10x smaller** error than MLMC

- We have  $\text{Cost} = O(\Delta^{-r})$  and estimated:

- $\hat{r}_{\text{NMC}} = 2.97$  ( $r_{\text{NMC}} = 3$ )

- $\hat{r}_{\text{NKQ}} = 1.9$  ( $r_{\text{NKQ}} = 2$ )

# Health economics

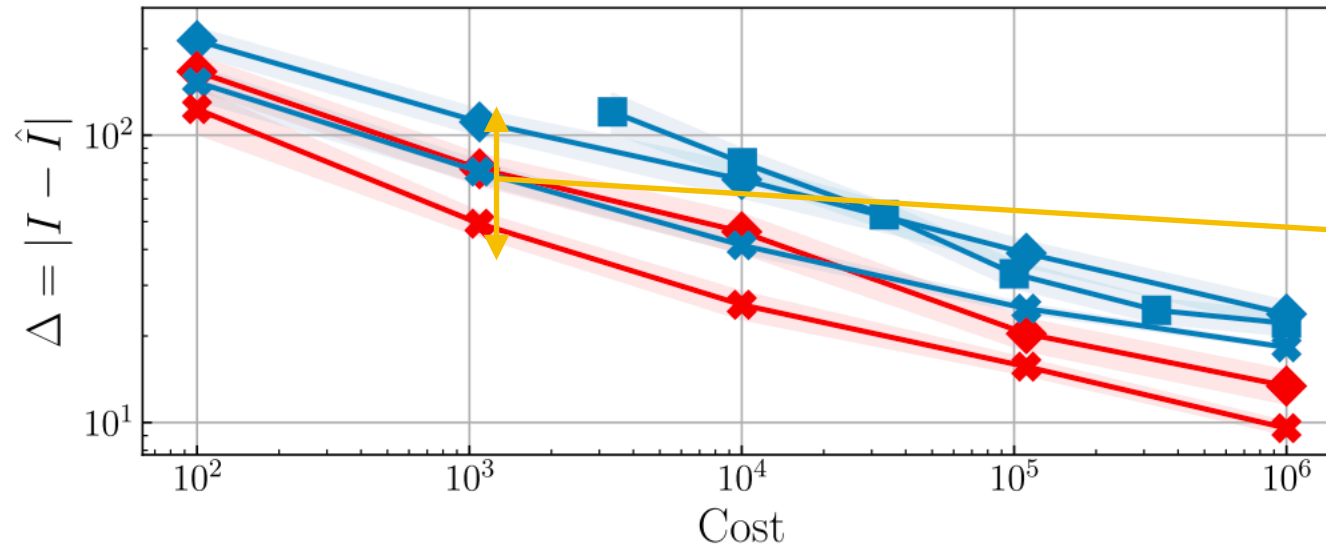


- **Problem:** Computing the expected value of partial perfect information. QoI for deciding whether we want to collect measurements on more variables from patients.

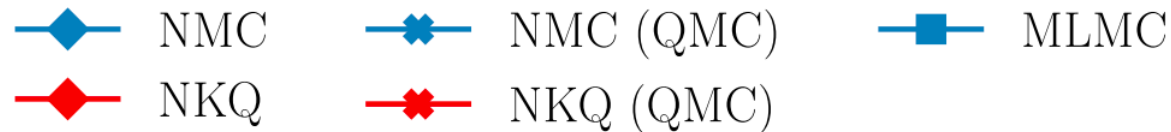
- We use  $s_{\mathcal{X}} = s_{\Theta} = \infty$ .

# Health economics

- **Problem:** Computing the expected value of partial perfect information. QoI for deciding whether we want to collect measurements on more variables from patients.

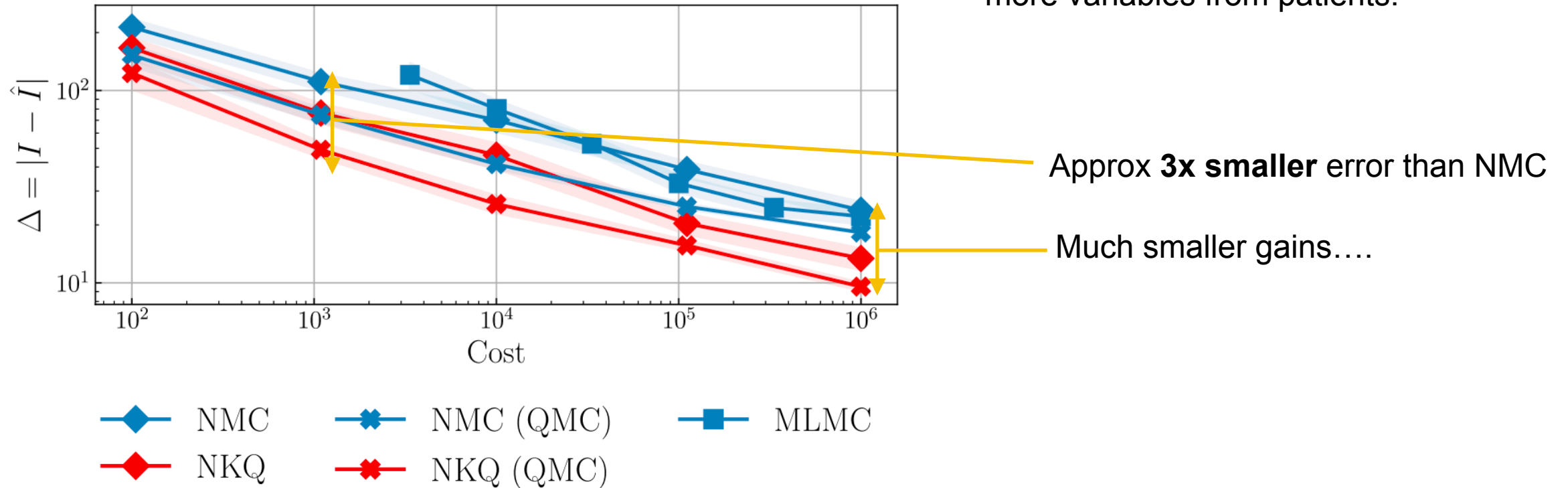


Approx **3x smaller** error than NMC



# Health economics

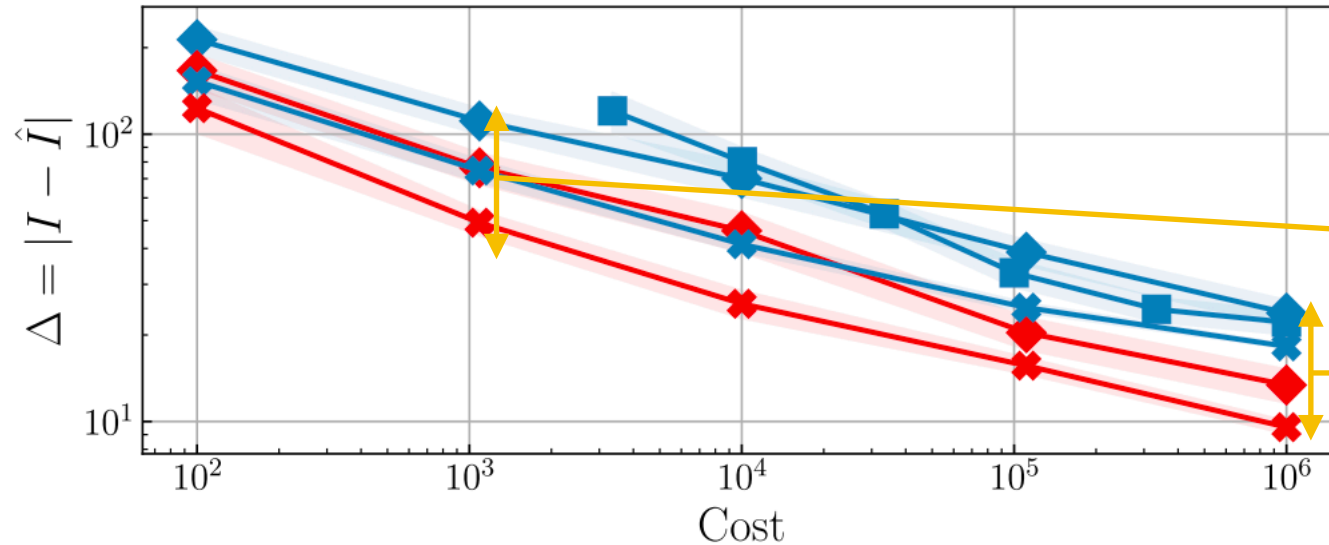
- **Problem:** Computing the expected value of partial perfect information. QoI for deciding whether we want to collect measurements on more variables from patients.





# Health economics

- **Problem:** Computing the expected value of partial perfect information. QoI for deciding whether we want to collect measurements on more variables from patients.



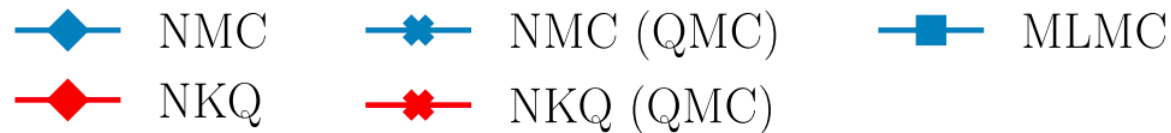
Approx **3x smaller** error than NMC

Much smaller gains....

## Why not as good?

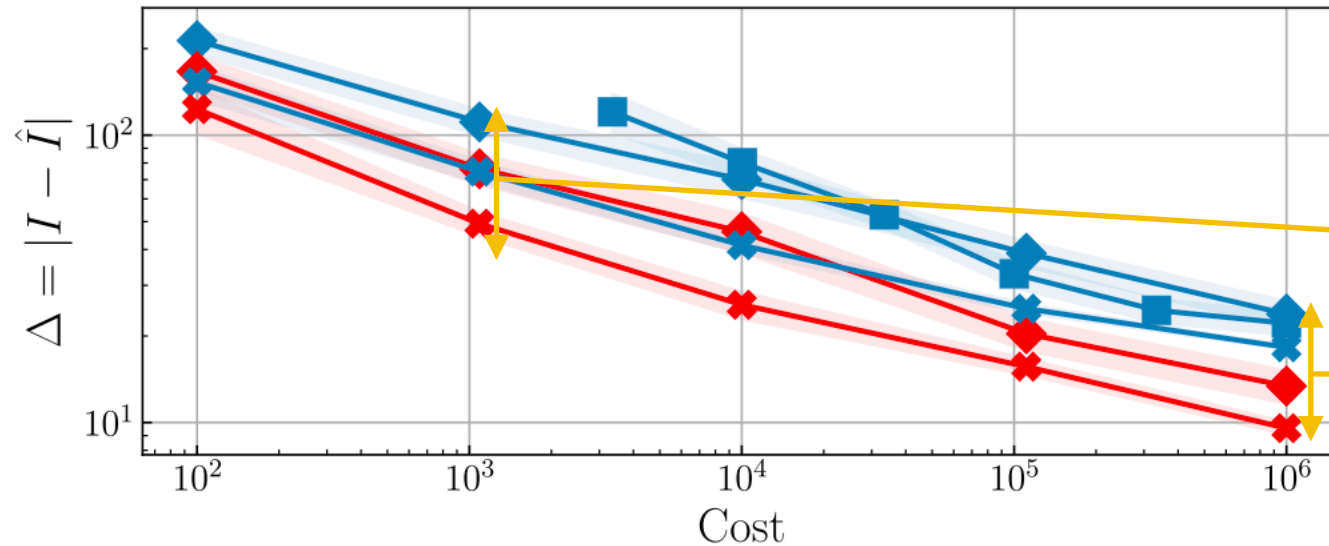
- The problem is high-dimensional:

$$d_x = 17$$



# Health economics

- **Problem:** Computing the expected value of partial perfect information. QoI for deciding whether we want to collect measurements on more variables from patients.



Approx **3x smaller** error than NMC

Much smaller gains....

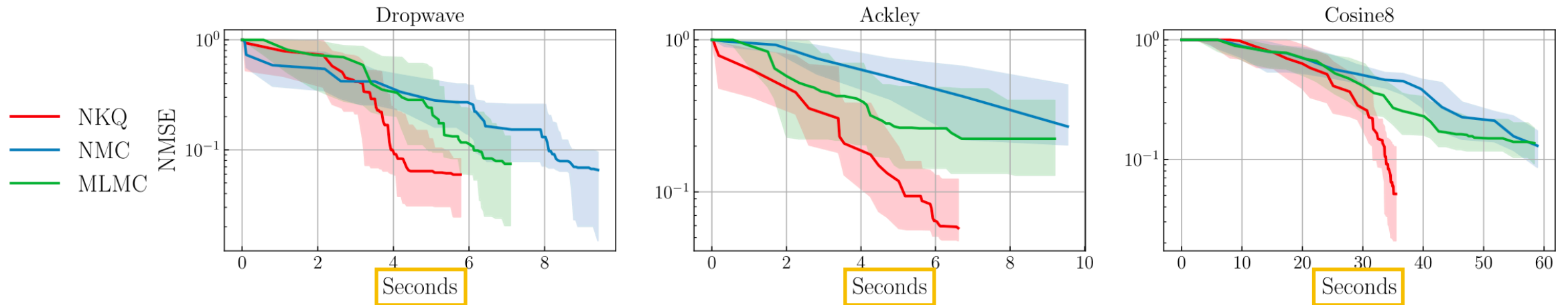
## Why not as good?

- The problem is high-dimensional:

$$d_x = 17$$

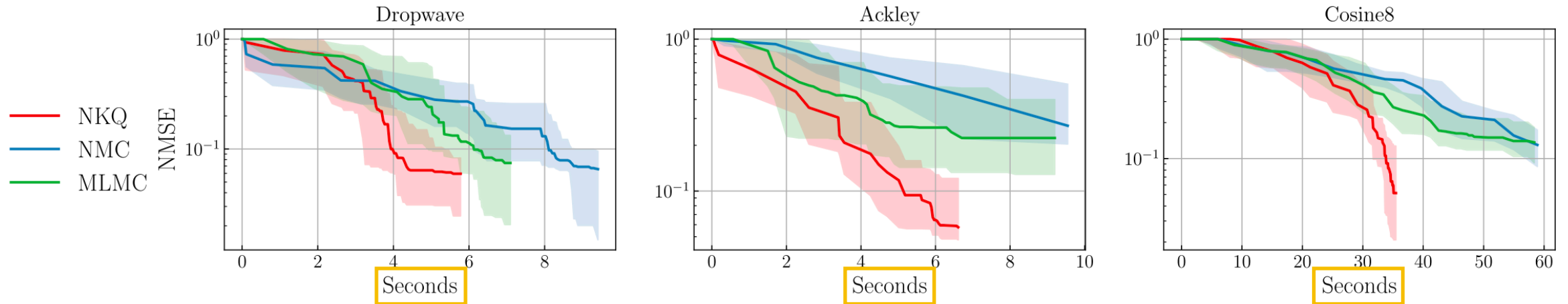
➡ We might still be very happy with 3x smaller error for small  $N$  &  $T$ !!

# Look-ahead Bayesian optimisation



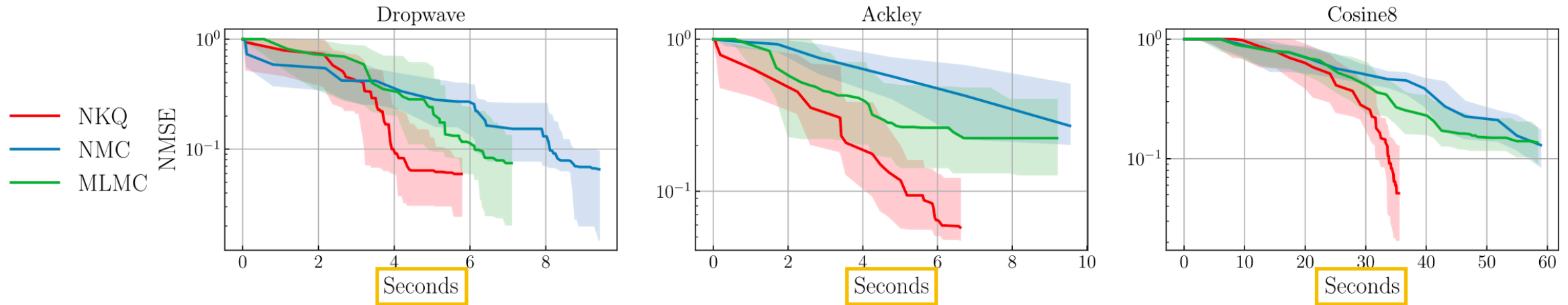
- **Problem:** Repeatedly computing and optimising acquisition function for 2-step look-ahead Bayesian optimisation. This requires compute a **very large number** of nested expectations.

# Look-ahead Bayesian optimisation



- **Problem:** Repeatedly computing and optimising acquisition function for 2-step look-ahead Bayesian optimisation. This requires compute a **very large number** of nested expectations.
- We have  $d_{\mathcal{X}} = d_{\Theta} = 2$ . We chose  $N = T = \Delta^{-2}$  for NMC, and  $N = T = \Delta^{-1}$  for NKQ.

# Look-ahead Bayesian optimisation



- **Problem:** Repeatedly computing and optimising acquisition function for 2-step look-ahead Bayesian optimisation. This requires compute a **very large number** of nested expectations.
- We have  $d_{\mathcal{X}} = d_{\Theta} = 2$ . We chose  $N = T = \Delta^{-2}$  for NMC, and  $N = T = \Delta^{-1}$  for NKQ.
- Cubic cost for NKQ only occurs once since the matrix inverses can be re-used (through change of variable trick).

# Summary and future work

- **Summary:** New estimator whose cost is orders of magnitude smaller than competitors when the integrands are smooth and dimensions not too high:

$$\text{Cost}(\hat{I}_{\text{NKQ}}) = \tilde{O}\left(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}-\frac{d_{\Theta}}{s_{\Theta}}}\right)$$

# Summary and future work

- **Summary:** New estimator whose cost is orders of magnitude smaller than competitors when the integrands are smooth and dimensions not too high:

$$\text{Cost}(\hat{I}_{\text{NKQ}}) = \tilde{O}\left(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}-\frac{d_{\Theta}}{s_{\Theta}}}\right)$$

- **Interesting extensions:**
  - Nested Bayesian quadrature? Useful for uncertainty quantification and active learning, but challenging propagation of uncertainty due to non-linearity of  $f$ .

# Summary and future work

- **Summary:** New estimator whose cost is orders of magnitude smaller than competitors when the integrands are smooth and dimensions not too high:

$$\text{Cost}(\hat{I}_{\text{NKQ}}) = \tilde{O}\left(\Delta^{-\frac{d_{\mathcal{X}}}{s_{\mathcal{X}}}-\frac{d_{\Theta}}{s_{\Theta}}}\right)$$

- **Interesting extensions:**
  - Nested Bayesian quadrature? Useful for uncertainty quantification and active learning, but challenging propagation of uncertainty due to non-linearity of  $f$ .
  - In-depth study of multilevel KQ approach?





**UCL**

# Any Questions?

Chen, Z., Naslidnyk, M., & Briol, F.-X. (2025). Nested expectations with kernel quadrature. *arXiv:2502.18284*.