

Nesting Probabilistic Programs

Tom Rainforth

With help from Rob Cornish, Hongseok Yang, Andrew Warrington, and Frank Wood

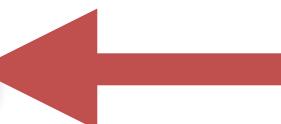


UNIVERSITY OF
OXFORD

Probabilistic programming presents and opportunity to nest queries

```
(defquery outer [D]
  (let [y (sample (beta 2 3))
        dist (conditional inner)
        z (sample (dist y D))])
  (* y z)))
```

```
(defquery inner [y D]
  (let [z (sample (gamma y 1))]
    (observe (normal y z) D)
    z))
```





Nested Expectations

$$\gamma_0 = \mathbb{E} [f_0 (y, \mathbb{E} [f_1 (y, z) | y])]$$

$$\gamma_0 = \mathbb{E}[f_0(y, \mathbb{E}[f_1(y, z, \mathbb{E}[f_2(y, z, u) | y, z] | y) | y])]$$

There are problems we cannot encode without nesting

- E.g. entropy of KL divergence of a marginal distribution:

$$H(p(y)) = -\mathbb{E} [\log (\mathbb{E} [p(y|\theta)|y])]$$

Agents reasoning about other agents



Aims

- When does nesting probabilistic programs no longer give a standard Bayesian model?
- What are the statistical implications?
- How can we construct valid inference engines?



Three Ways We Might Nest

- Nested conditioning
- Nested inference
- Estimates as first class variables

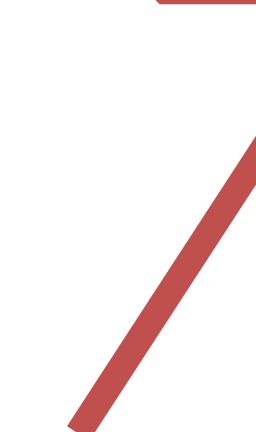


Nested Conditioning

```
(defquery outer [D]
  (let [y (sample (beta 2 3))])
    (observe (nest inner [y D] :smc 100) nil)
    y))
```

Formalizing Nested Conditioning

$$\pi_o(y) = \boxed{\psi(y)} \boxed{p_i(y)} \approx \psi(y) \hat{p}_i(y)$$




Density terms from outer

Marginal likelihood for inner



Formalizing Nested Conditioning

$$\mathbb{E}_{q(y)} \left[\frac{\psi(y) \hat{p}_i(y)}{q(y)} \right] = \mathbb{E}_{q(y)} \left[\frac{\pi_o(y)}{q(y)} \right]$$

Nested Inference

```
(defquery inner [y D]
  (let [z (sample (gamma y 1)) ]
    (observe (normal y z) D)
    z))
```

```
(defquery outer [D]
  (let [y (sample (beta 2 3))
        dist (conditional inner)
        z (sample (dist y D)) ]
    (* y z)))
```



Formalizing Nested Inference

$$\pi_o(y, z) = \psi(y, z)p_i(z|y)$$

$$p_i(z|y) = \frac{\pi_i(y, z)}{\int \pi_i(y, z') dz'}$$

$$\pi_o(y, z) = \frac{\psi(y, z)\pi_i(y, z)}{\int \pi_i(y, z') dz'}$$

Formalizing Nested Inference

$$\pi_o(y, z) \approx \psi(y, z) \hat{p}_i(z|y)$$

$$\mathbb{E}_{q(y, z)} \left[\frac{\pi_o(y, z)}{q(y, z)} \right] = \mathbb{E}_{q(y, z)} \left[\frac{\psi(y, z) \pi_i(y, z)}{q(y, z) \mathbb{E}_{z' \sim q(z|y)} [\pi_i(y, z') / q(z'|y)]} \right]$$



Estimates as First Class Variables

```
(defm prior [] (normal 0 1))
(defm lik [theta d] (normal theta d))

(defquery inner-q [y d]
  (let [theta (sample (prior))]
    (observe (lik theta d) y)))

(defn inner-E [y d M]
  (->> (doquery :importance
    inner-q [y d])
    (take M)
    log-marginal))

(with-primitive-procedures [inner-E]

(defquery outer-q [d M]
  (let [theta (sample (prior))
        y (sample (lik theta d))
        log-lik (observe*
                  (lik theta d) y)
        log-marg (inner-E y d M)]
    (- log-lik log-marg)))))

(defn outer-E [d M N]
  (->> (doquery :importance
    outer-q [d M])
    (take N)
    collect-results
    empirical-mean))
```



Estimates as First Class Variables

$$\pi_0(y) = \psi(y, z(y)) \approx \psi(y, \hat{z}(y))$$

$$\mathbb{E}_{q(y)} \left[\frac{\pi_0(y)}{q(y)} \right] = \mathbb{E}_{q(y)} \left[\frac{\psi(y, \mathbb{E}[\hat{z}(y)|y])}{q(y)} \right]$$

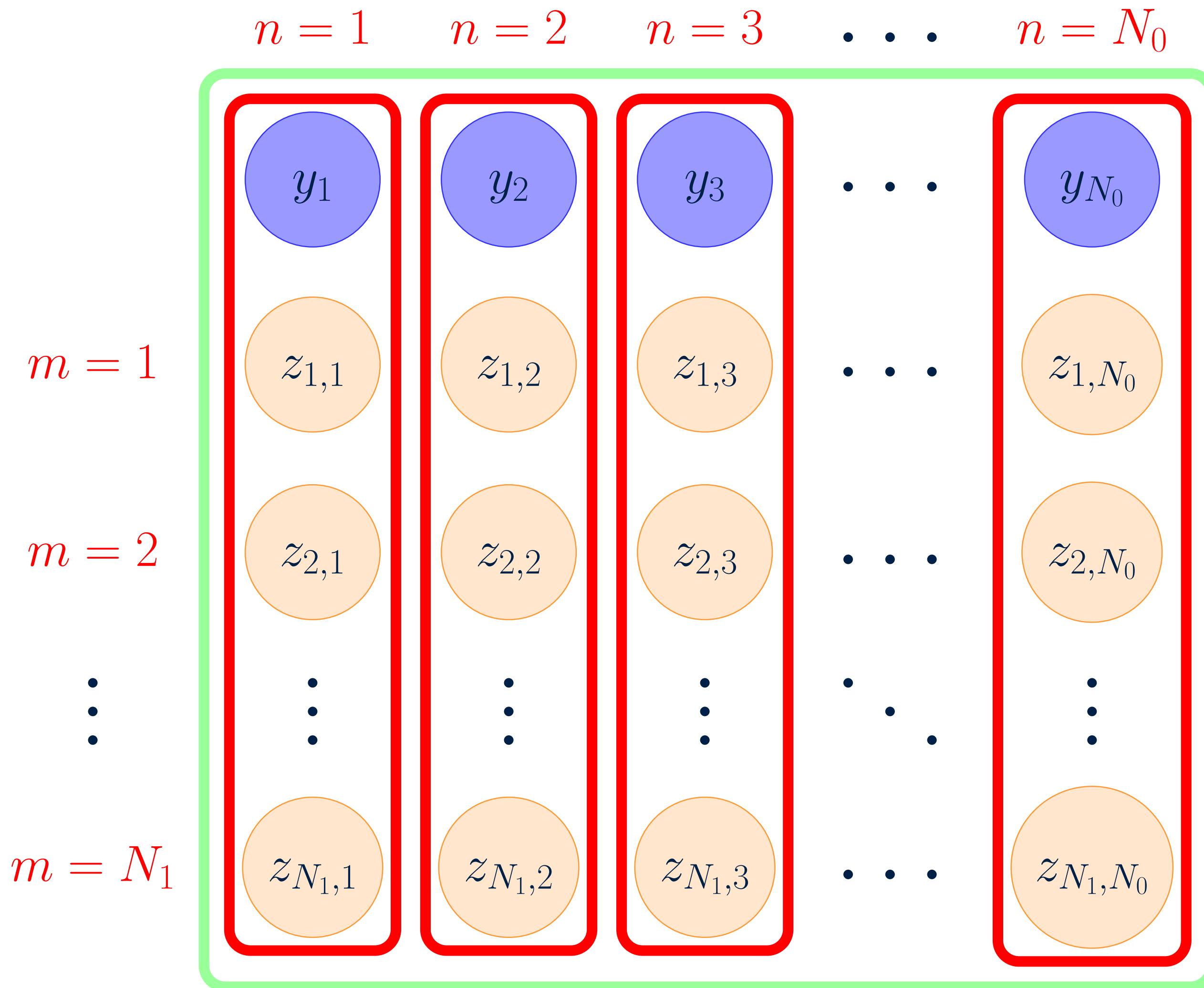


Nested Monte Carlo

$$I_0 = \frac{1}{N_0} \sum_{n=1}^{N_0} f_0(y_n, I_1^n(y_n)) \quad \text{where} \quad y_n \sim p(y)$$

$$I_1^n(y_n) = \frac{1}{N_1} \sum_{m=1}^{N_1} f_1(y_n, z_{m,n}) \quad \text{where} \quad z_{m,n} \sim p(z|y = y_n)$$

[Rainforth, Cornish, Yang, Warrington, and Wood. On the Opportunities and Pitfalls of Nesting Monte Carlo Estimators. ICML 2018]



Convergence Rate

- If each f_k is continuously differentiable, the MSE converges at a rate

$$O\left(\frac{1}{N_0} + \left(\sum_{k=1}^D \frac{1}{N_k}\right)^2\right)$$



Variance Bias Squared

- We need each $N_k \rightarrow \infty$ for convergence



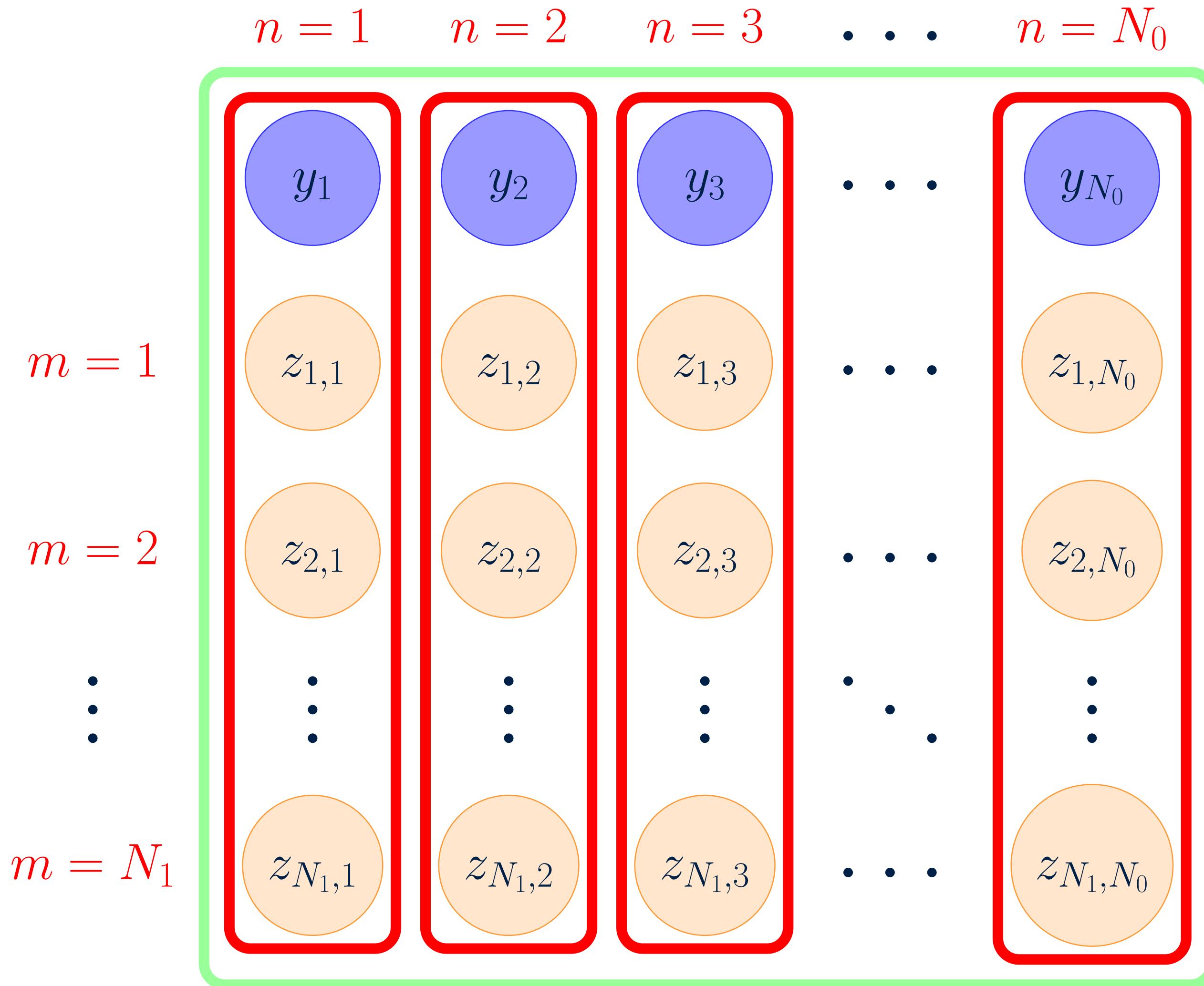
Optimal Setup

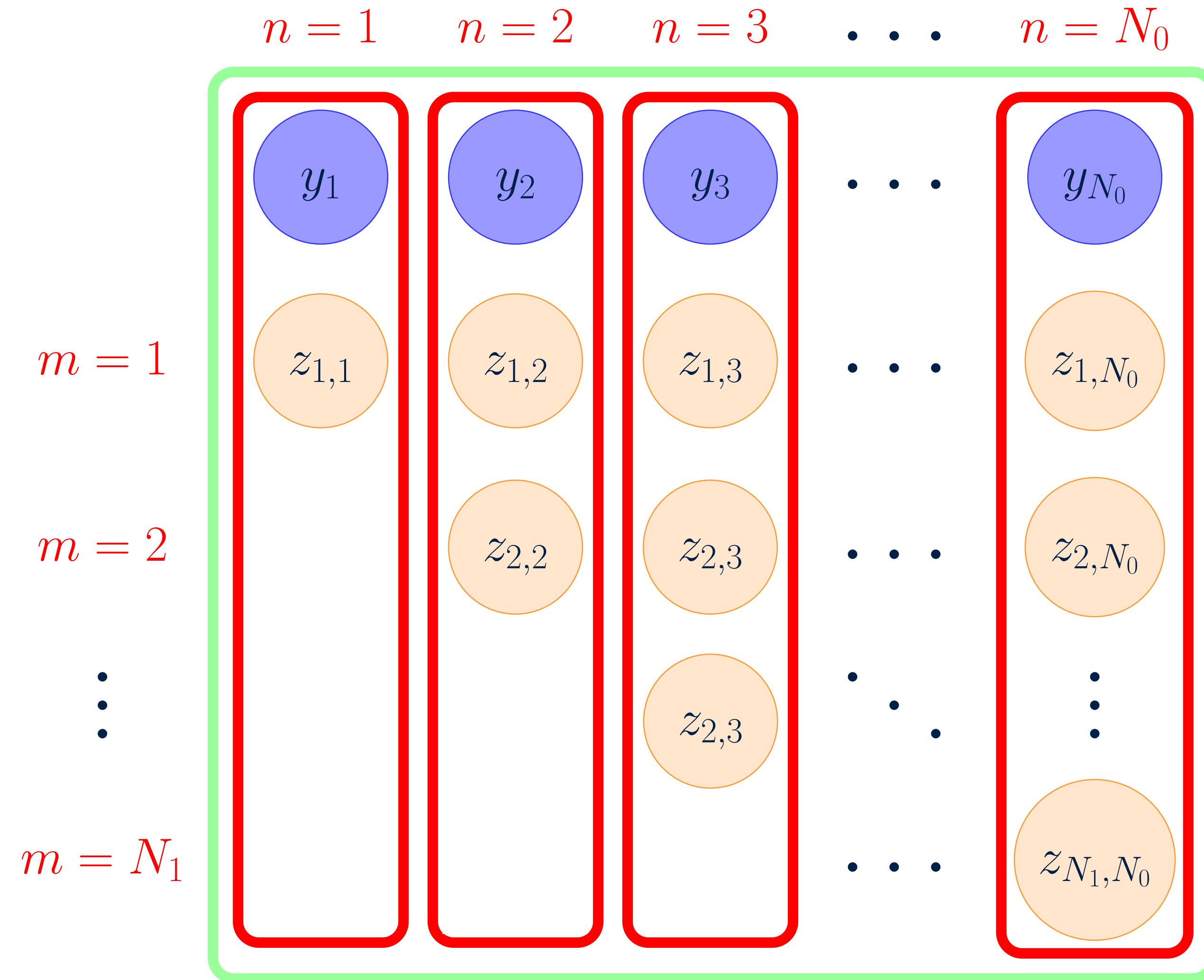
- The asymptotic mean squared error is minimised when

$$\sqrt{N_0} \propto N_1 \propto N_2 \propto \dots \propto N_D$$

- This yields an overall convergence rate of $O\left(\frac{1}{T^{\frac{2}{2+D}}}\right)$

$$D = 0, O\left(\frac{1}{T}\right); \quad D = 1, O\left(\frac{1}{T^{\frac{2}{3}}}\right); \quad D = 2, O\left(\frac{1}{T^{\frac{1}{2}}}\right)$$





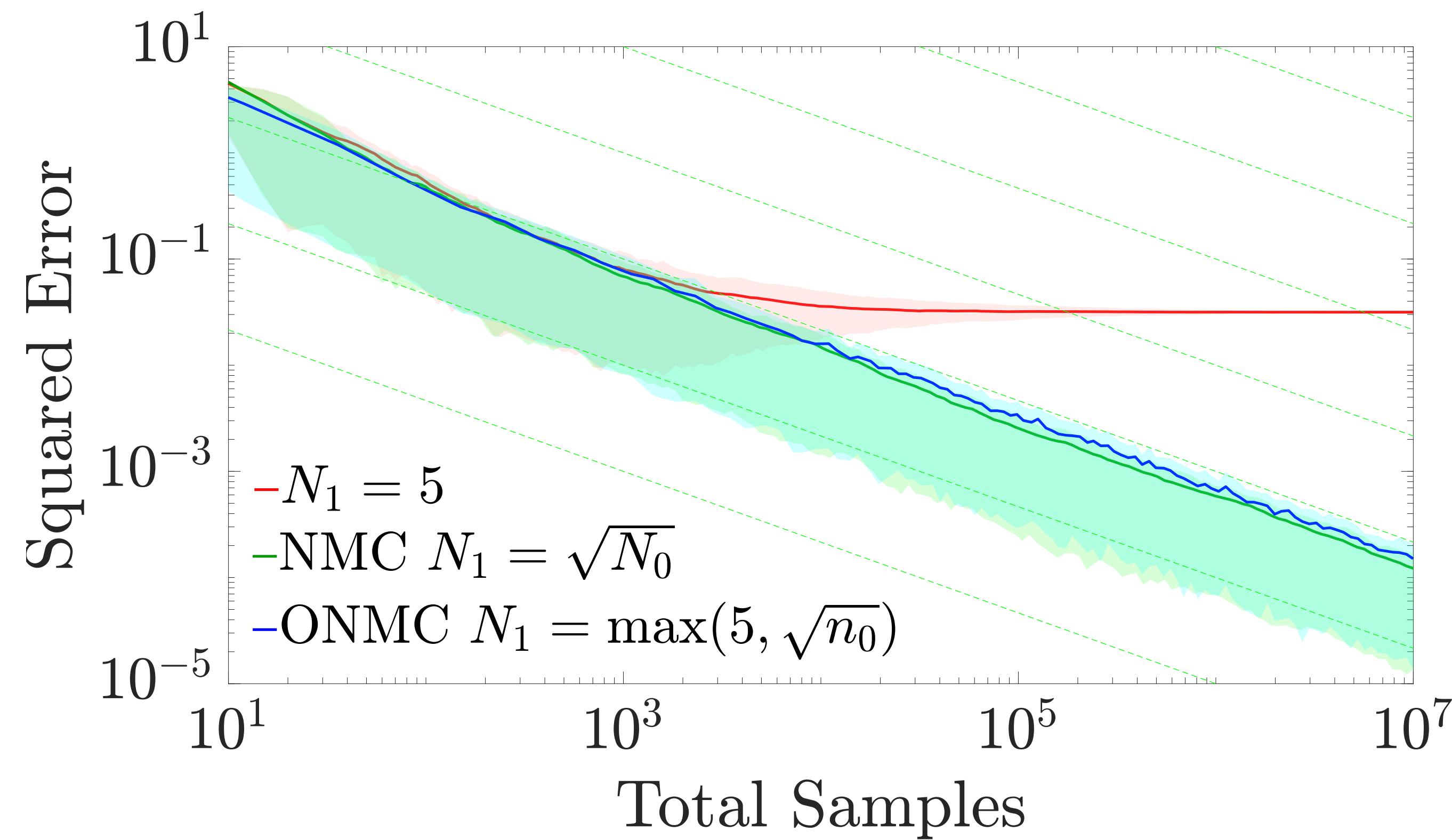


Online NMC

$$\text{Var}[\text{ONMC}] \rightarrow c \text{Var}[\text{NMC}] \quad \text{where } c < 1$$

$$\text{Bias}[\text{ONMC}] \rightarrow d \text{Bias}[\text{NMC}] \quad \text{where } d < 2$$

Online NMC



Take Homes

- Nesting probabilistic programs allows encoding of models we could not otherwise express
 - Allows encoding of nested expectations
- We need to be careful in the design of our inference engines
- Online NMC gives an easy way of making sure consistency is maintained



Thanks!



Rob Cornish



Hongseok Yang



Andrew Warrington



Frank Wood

[Rainforth, Cornish, Yang, Warrington, and Wood. On the Opportunities and Pitfalls of Nesting Monte Carlo Estimators. ICML 2018]

[Rainforth. Nesting Probabilistic Programs. UAI 2018]

Contact: rainforth@stats.ox.ac.uk