

- ▶ Basic Bayesian statistics
 - ▶ Prior, Observations, Posterior, and Bayes' formula
 - ▶ Computational approaches
 - ▶ Priors
 - ▶ Example: Beta&Binomial
 - ▶ Credible intervals vs Confidence intervals
- ▶ Integration methods
 - ▶ Deterministic integration
 - ▶ Monte Carlo integration

Bayesian statistics

Prior, Observations, Posterior; concepts

In a basic Bayesian model, the main difference to frequentistic modelling is that the unknown parameters are now treated as *unobserved random variables* instead of unknown constants.

- ▶ The distribution for the parameters is called the *prior distribution*.
- ▶ The observation distribution conditionally on the parameters is often referred to as the observation *likelihood*, although this terminology is not universally accepted as proper.
- ▶ The *conditional distribution* for the parameters, given the observations, is called the *posterior distribution*

Prior, Observations, Posterior; mathematical probability theory

In mathematical terms, with parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)$ and observations $\mathbf{y} = (y_1, \dots, y_n)$:

- ▶ Prior density: $\boldsymbol{\theta} \sim p_{\boldsymbol{\theta}}(\boldsymbol{\theta})$, often abbreviated to just $p(\boldsymbol{\theta})$. If the prior distributions have any (fixed and known) parameters, those are often referred to as *hyper-parameters*.
- ▶ Conditional observation likelihood/density/probability mass function: $(\mathbf{y}|\boldsymbol{\theta}) \sim p_{\mathbf{y}|\boldsymbol{\theta}}(\mathbf{y})$, usually written $p(\mathbf{y}|\boldsymbol{\theta})$.
- ▶ Posterior density: $(\boldsymbol{\theta}|\mathbf{y}) \sim p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\boldsymbol{\theta}, \mathbf{y})}{p(\mathbf{y})} = \frac{p(\boldsymbol{\theta})p(\mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{y})}$ Bayes' formula
- ▶ *Marginal* observation density/probability mass function:
$$p(\mathbf{y}) = \int p(\boldsymbol{\theta}, \mathbf{y}) d\boldsymbol{\theta} = \int p(\boldsymbol{\theta})p(\mathbf{y}|\boldsymbol{\theta}) d\boldsymbol{\theta}$$

For discrete variables or parameters, the integrals are replaced by sums.

Computational approaches

The Bayesian statistics framework is appealing in that (almost) everything flows from general probability theory, once a probabilistic model has been formulated. The difficulty lies in how to actually calculate the involved densities, expectations, variances, and probabilities.

- ▶ Analytical derivation; some results are easy to derive by hand (see example on later slide)
- ▶ Approximate deterministic numerical integration; deterministic integration schemes, Laplace approximation
- ▶ Basic simulation and integration; Monte Carlo integration, importance sampling
- ▶ Advanced simulation and integration; Markov chain Monte Carlo (MCMC) simulation (beyond the scope of this course)

Priors

One of the challenges in Bayesian modelling is how to specify realistic priors.

- ▶ One approach is to elicit information from application area experts and encode their *subjective* information as probability distributions.
- ▶ Other approaches include so called *objective* priors, *vague* priors, and *flat* priors (often incorrectly called "uninformative" priors) with the aim of regularising otherwise unstable results but without imposing too strict limitations.
- ▶ Before the proliferation of computational techniques in the 1990's, and indeed still today, most priors were chosen to be *conjugate* with respect to the observation distribution (the prior and posterior belong to the same distribution family), to allow analytical solutions. Modern computational methods has removed much of the motivation for such priors; they can simplify calculations, but need not be used if more realistic priors are available.

Example: Binomial observations with Beta prior

Consider the following Bayesian model (version with a single observation seen in Statistical Methodology lectures) for Binomial observations $\mathbf{y} = (y_1, \dots, y_K)$ with common probability parameter ϕ :

- ▶ Parameter prior distribution: $\phi \sim \text{Beta}(a, b)$ for some *hyper-parameters* $a, b > 0$.
- ▶ Observations: $(y_k|\phi) \sim \text{Bin}(N_k, \phi)$, independent over $k = 1, \dots, K$, known N_k .
- ▶ Posterior distribution for ϕ : $\text{Beta}\left[a + \sum_{k=1}^K y_k, b + \sum_{k=1}^K (N_k - y_k)\right]$

Proof: Moving all factors that do not depend on ϕ into a normalisation constant, we get

$$\begin{aligned} p(\phi|\mathbf{y}) &= \frac{p(\phi)p(\mathbf{y}|\phi)}{p(\mathbf{y})} \propto p(\phi) \prod_{k=1}^K p(y_k|\phi) = \frac{\phi^{a-1}(1-\phi)^{b-1}}{B(a,b)} \prod_{k=1}^K \binom{N_k}{y_k} \phi^{y_k} (1-\phi)^{N_k-y_k} \\ &\propto \phi^{a-1+\sum_k y_k} (1-\phi)^{b-1+\sum_k (N_k-y_k)}. \end{aligned}$$

This takes the same form as a Beta distribution density (this shows that Beta is a conjugate prior for the probability parameter in a Binomial distribution), so the result follows by identifying the parameters.

Bayesian credible intervals

A Bayesian $(1 - \alpha) \cdot 100$ percent credible interval $\text{CI}_\theta = (L, U)$ for a parameter θ is an interval computed from the posterior distribution, such that $P(L < \theta < U | \mathbf{y}) \geq 1 - \alpha$.

- ▶ Often, L and U are chosen so that the tail probabilities outside the interval are both $= \alpha/2$.
- ▶ Another option is to choose the smallest set or interval with the required probability. This is achieved by finding a value c such that the set $\text{CI}_\theta = \{\theta; p(\theta|\mathbf{y}) > c\}$ fulfills $P(\theta \in \text{CI}_\theta | \mathbf{y}) \geq 1 - \alpha$. The resulting CI_θ is called a *highest posterior density* (HPD) region, and is well defined also for multi-dimensional parameter vectors.
- ▶ When c is chosen as large as possible, the HPD is the smallest credible region possible for the given model parameterisation. However, due to how probability densities are changed under transformation of the random variable, HPD region constructions are not invariant under reparameterisation.

Credible region vs. Confidence region

Posterior credible region (Bayesian concept)

A *level $1 - \alpha$ posterior credible region* for θ , $C_\theta(\mathbf{y})$, is a set (usually an interval) such that

$$\mathbb{P}_{\theta \sim p(\theta|\mathbf{y})} (\theta \in C_\theta(\mathbf{y})) \geq 1 - \alpha$$

for the fixed set of observations \mathbf{y} . The probability is a direct statement about our posterior beliefs about the random variable θ .

Confidence region (Frequentist concept)

A *level $1 - \alpha$ confidence region* procedure for θ , $C_\theta(\mathbf{y})$, generates sets (usually intervals) in such a way that

$$\mathbb{P}_{\mathbf{y} \sim p(\mathbf{y}; \theta)} (\theta \in C_\theta(\mathbf{y})) \geq 1 - \alpha$$

for every possible θ (or at least for the true value). The probability statement concerns the confidence region construction *procedure* under repeated experiments; the observations \mathbf{y} are random.

Numerical integration

Numerical integration

Non-statistical and statistical problems involving integrals:

- ▶ Integrate $f(\cdot)$ over some interval/set/domain Ω :

$$I = \int_{\Omega} f(\mathbf{x}) \, d\mathbf{x}$$

- ▶ Compute a *marginal density* $p_Y(y)$ from a *joint density* $p_{X,Y}(x,y)$, $x \in \Omega$:

$$p_Y(y) = \int_{\Omega} p_{X,Y}(x,y) \, dx$$

- ▶ For some function $\phi(\mathbf{x})$ and density $p_X(\mathbf{x})$, $\mathbf{x} \in \Omega$, compute the expectation

$$\mathbb{E}_X[\phi(\mathbf{x})] = \int_{\Omega} \phi(\mathbf{x}) p_X(\mathbf{x}) \, d\mathbf{x}$$

- ▶ Challenges: High-dimensional Ω and computationally expensive integrands

Deterministic integration methods

- General idea: Find points and weights (\mathbf{x}_k, w_k) such that

$$I = \int_{\Omega} f(\mathbf{x}) \, d\mathbf{x} \approx \sum_k f(\mathbf{x}_k) w_k = \hat{I}$$

- Midpoint rule for $\Omega = (a, b)$: $x_k = a + (k - \frac{1}{2}) \frac{b-a}{N}$, $w_k = \frac{b-a}{N}$, for $k = 1, \dots, N$.
- Trapezoidal rule
- Simpson's rule
- Gaussian quadrature: For some special function $\psi(x)$, there is a specific set $\{(x_k, w_k)\}_N$ such that $\int g(x) \psi(x) \, dx = \sum_{k=1}^N g(x_k) w_k$ for any polynomial of degree $\leq 2N - 1$. Use

$$\hat{I} = \sum_{k=1}^N \frac{f(x_k)}{\psi(x_k)} w_k; \quad \text{good approximation of } I \text{ if } f(x)/\psi(x) \text{ is close to a polynomial.}$$

- Laplace approximation for $f(\mathbf{x}) \geq 0$ on $\Omega = \mathbb{R}^d$: Scale a Gaussian density to match the amplitude and shape at the mode of the integrand:

$$\hat{\mathbf{x}} = \underset{\mathbf{x} \in \Omega}{\operatorname{argmax}} f(\mathbf{x}), \quad \widehat{\mathbf{H}} = -\nabla^2 \log f(\mathbf{x})|_{\mathbf{x}=\hat{\mathbf{x}}}, \quad \hat{w} = \frac{(2\pi)^{d/2}}{(\det \widehat{\mathbf{H}})^{1/2}}, \quad \hat{I} = f(\hat{\mathbf{x}}) \hat{w}$$

Monte Carlo integration

- We know that if we have a sample $\mathbf{x}_1, \dots, \mathbf{x}_N$ from a distribution with density $p_X(\mathbf{x})$ we can estimate the expectation $E_{\mathbf{x} \sim p_X}[\phi(\mathbf{x})] = \int \phi(\mathbf{x}) p_X(\mathbf{x}) d\mathbf{x}$ with the average of $\phi(\mathbf{x}_k)$:

$$\hat{E}_{\mathbf{x} \sim p_X}[\phi(\mathbf{x})] = \frac{1}{N} \sum_{k=1}^N \phi(\mathbf{x}_k)$$

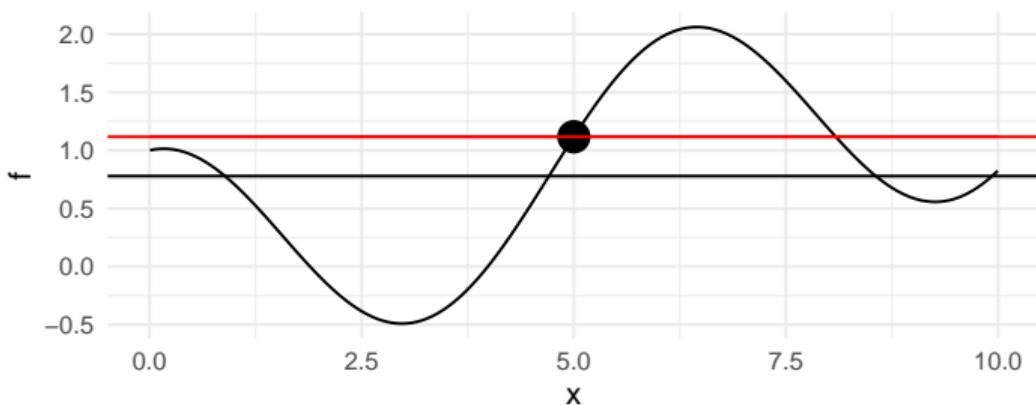
- For bounded Ω we can write a Monte Carlo integration scheme for

$$I = \int_{\Omega} f(\mathbf{x}) d\mathbf{x} = \int_{\Omega} |\Omega| f(\mathbf{x}) \frac{1}{|\Omega|} d\mathbf{x} \approx \frac{|\Omega|}{N} \sum_{k=1}^N f(\mathbf{x}_k), \quad \mathbf{x}_1, \dots, \mathbf{x}_N \sim \text{Unif}(\Omega)$$

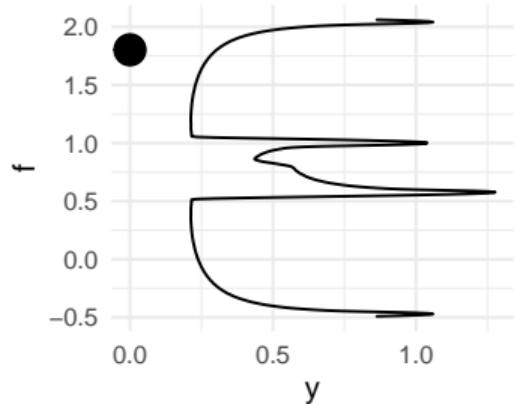
- Plain MC estimates are unbiased, with variance

$$\text{Var}_{\{\mathbf{x}_k \sim \text{Unif}(\Omega)\}} \left[\frac{|\Omega|}{N} \sum_{k=1}^N f(\mathbf{x}_k) \right] = \frac{|\Omega|^2}{N^2} N \text{Var}_{\mathbf{x} \sim \text{Unif}(\Omega)} [f(\mathbf{x})] \propto N^{-1}$$

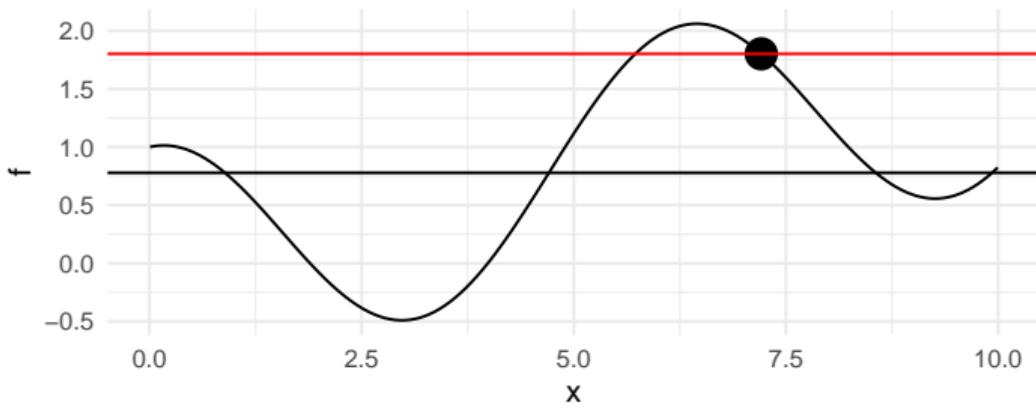
Midpoint, $N = 1$



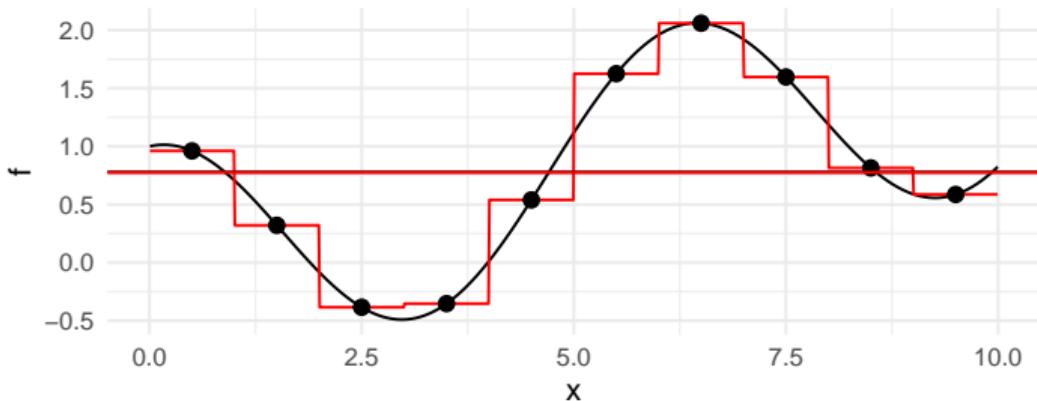
Monte Carlo



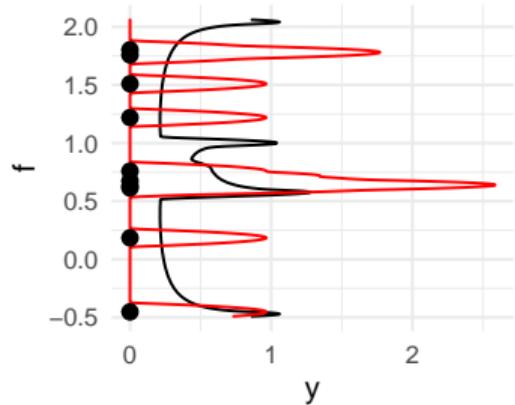
Monte Carlo, $N = 1$



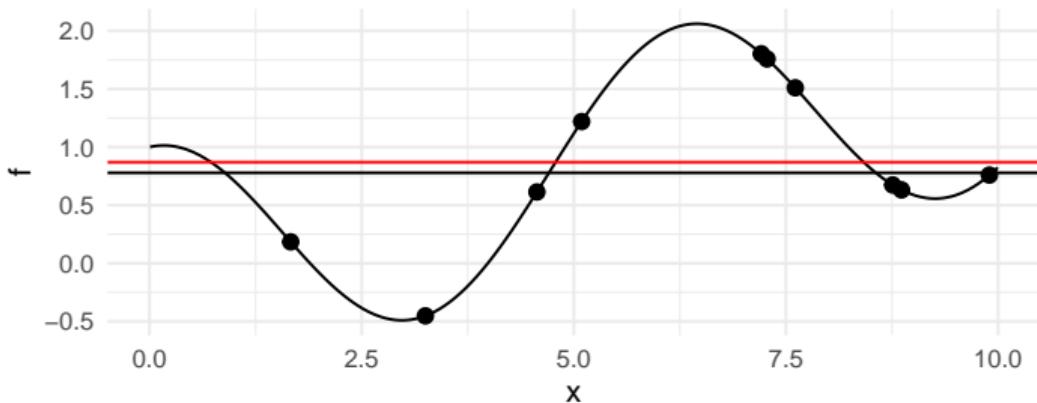
Midpoint, $N = 10$



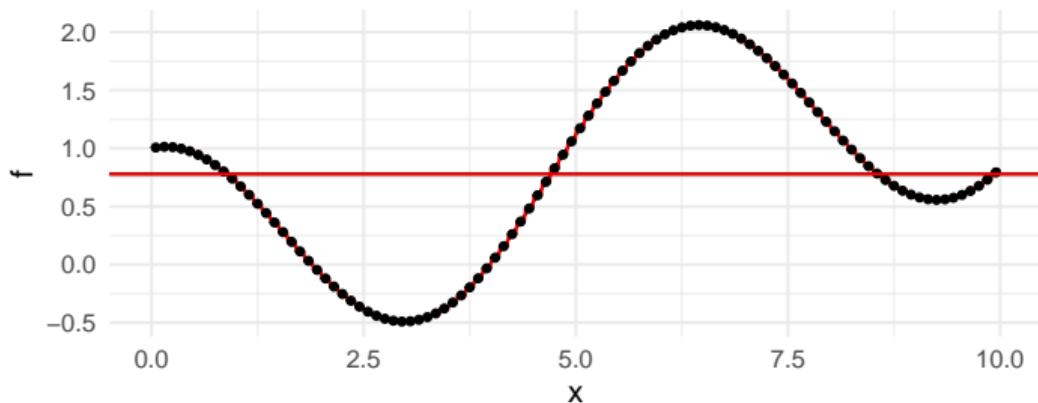
Monte Carlo



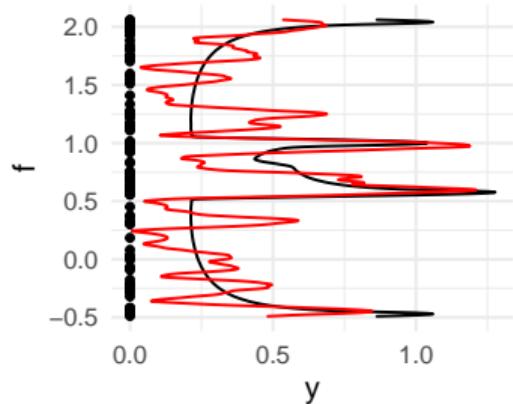
Monte Carlo, $N = 10$



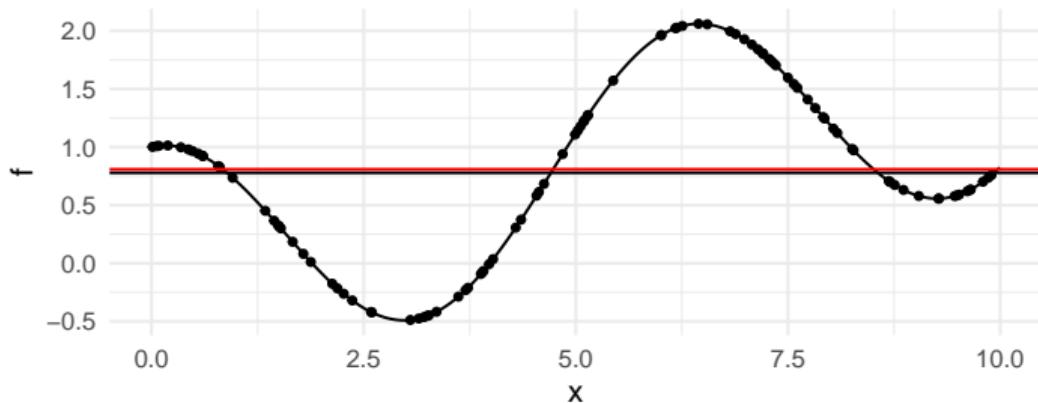
Midpoint, $N = 100$



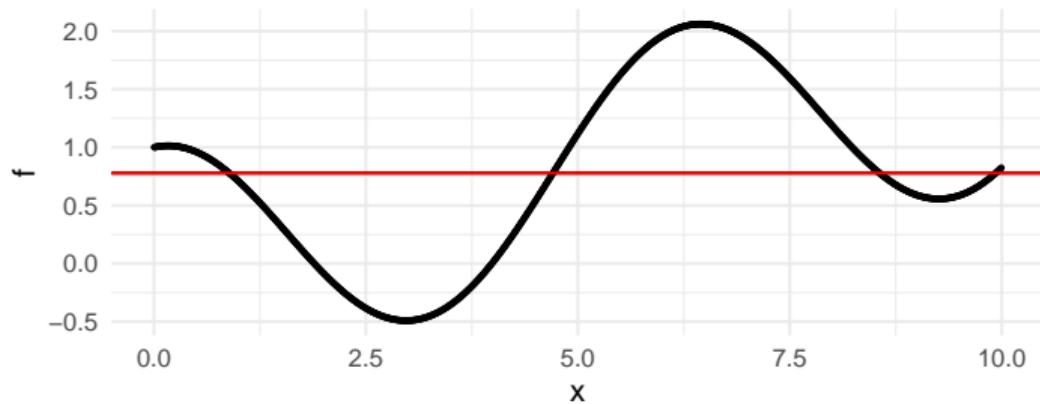
Monte Carlo



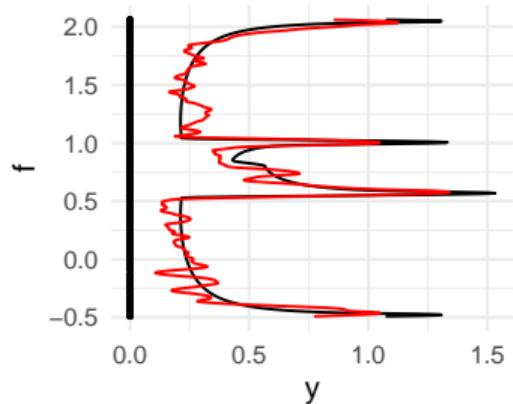
Monte Carlo, $N = 100$



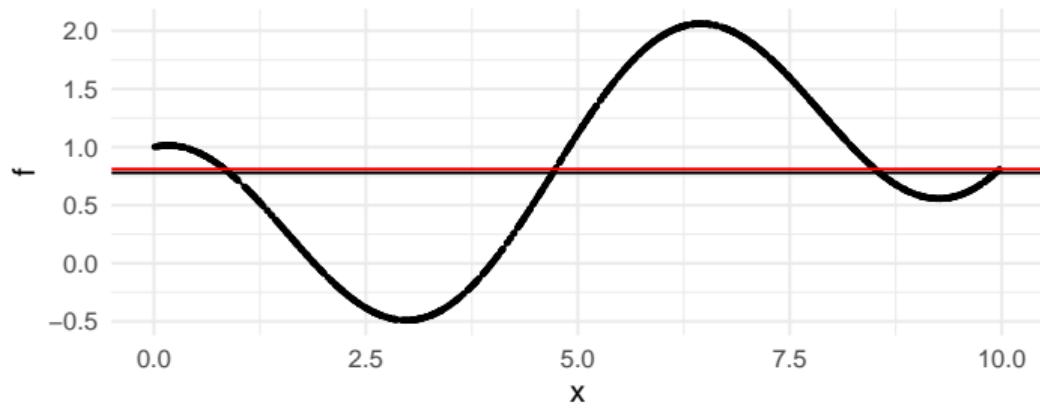
Midpoint, $N = 1000$



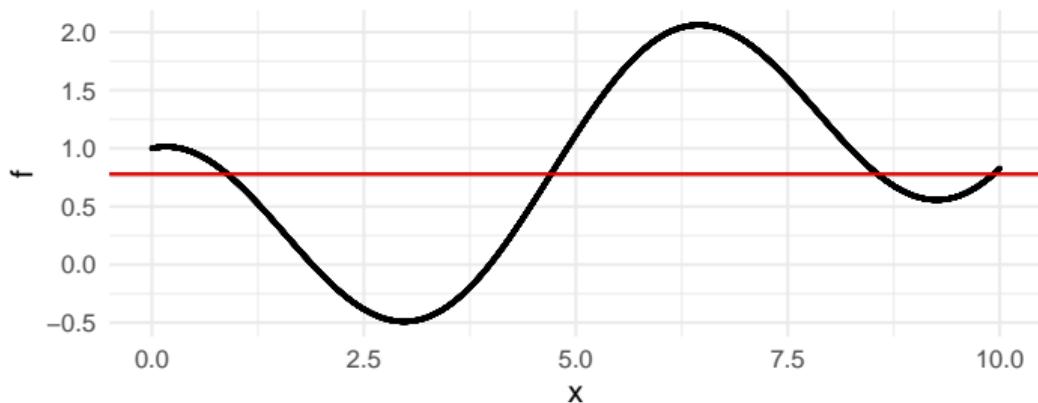
Monte Carlo



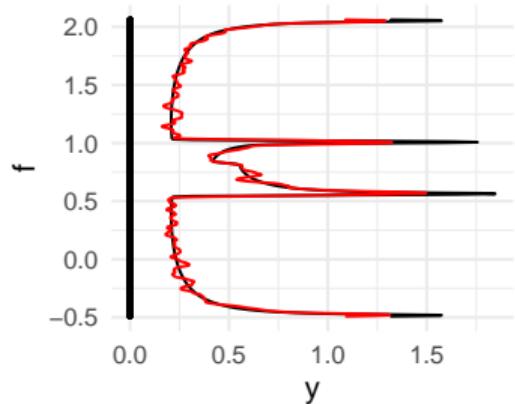
Monte Carlo, $N = 1000$



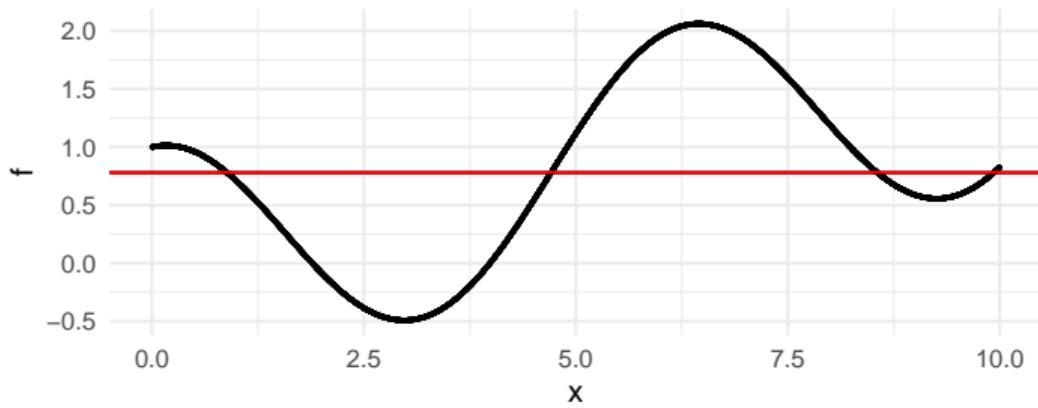
Midpoint, $N = 10000$



Monte Carlo



Monte Carlo, $N = 10000$



Importance sampling

- We can often reduce the MC variance by sampling with a density $p_Z(\cdot)$ that is similar to $\phi(\mathbf{x})p_X(\mathbf{x})$ or $f(\mathbf{x})$, a technique called *Importance sampling*:

$$I = \int f(\mathbf{x}) d\mathbf{x} = \int \frac{f(\mathbf{x})}{p_Z(\mathbf{x})} p_Z(\mathbf{x}) d\mathbf{x} = \mathbb{E}_{\mathbf{x} \sim p_Z} \left[\frac{f(\mathbf{x})}{p_Z(\mathbf{x})} \right] \approx \frac{1}{N} \sum_{k=1}^N \frac{f(\mathbf{x}_k)}{p_Z(\mathbf{x}_k)} = \hat{I}$$

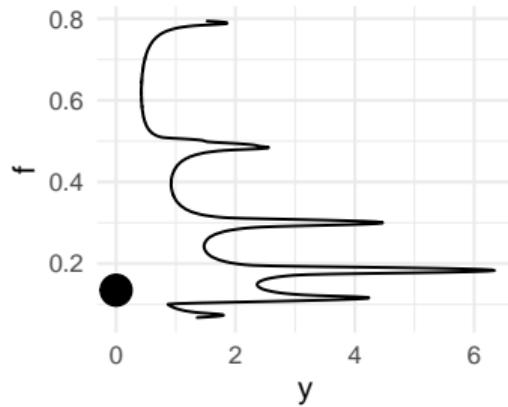
where \mathbf{x}_k , $k = 1, \dots, N$ are sampled from the $p_Z(\cdot)$ density.

- The variance is still $\propto N^{-1}$, but with a potentially much smaller constant:

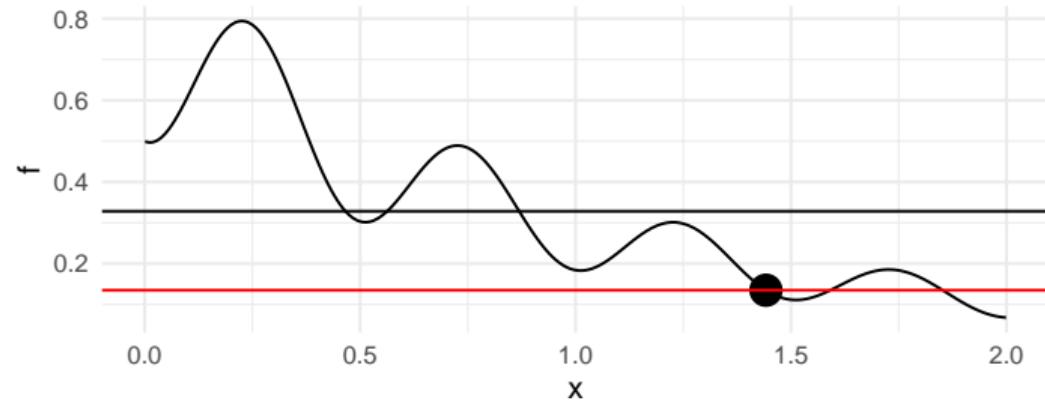
$$\text{Var}_{\{\mathbf{x}_k \sim p_Z\}} [\hat{I}] = \text{Var}_{\{\mathbf{x}_k \sim p_Z\}} \left[\frac{1}{N} \sum_{k=1}^N \frac{f(\mathbf{x}_k)}{p_Z(\mathbf{x}_k)} \right] = \frac{1}{N} \text{Var}_{\mathbf{x} \sim p_Z} \left[\frac{f(\mathbf{x})}{p_Z(\mathbf{x})} \right]$$

- Note: If we were able to choose $p_Z(\mathbf{x}) \propto f(\mathbf{x})$, the variance would be zero!

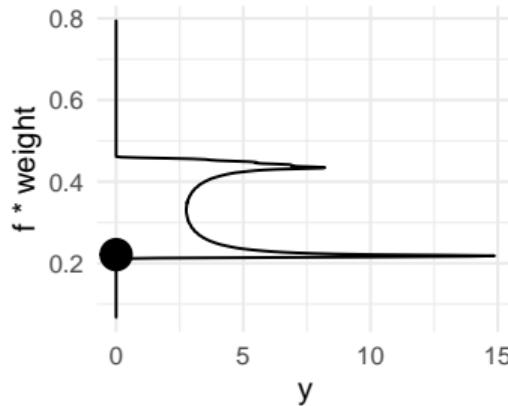
Monte Carlo



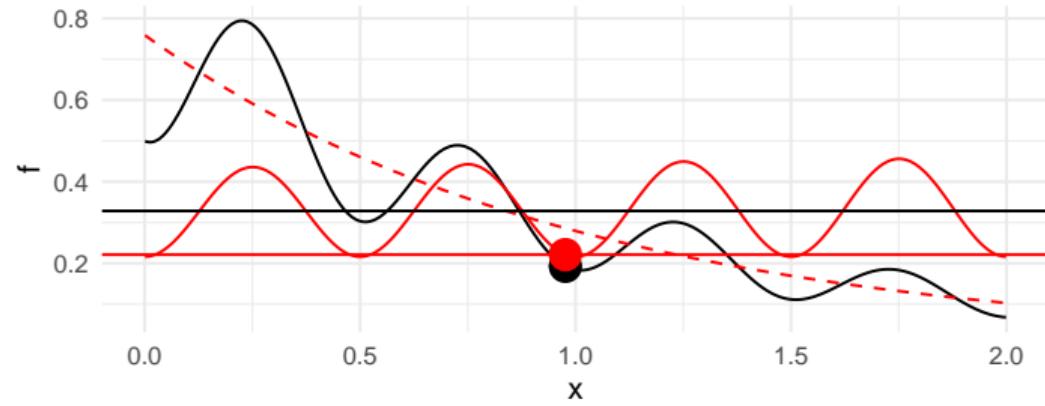
Monte Carlo, $N = 1$



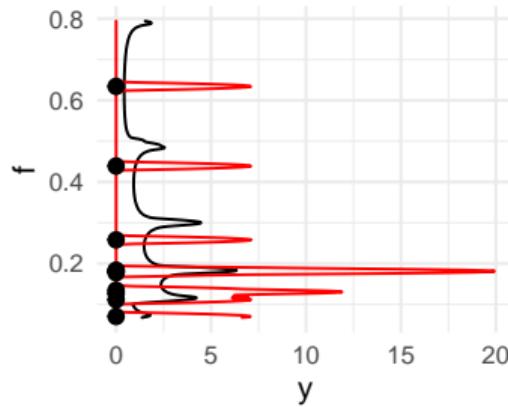
MC Importance



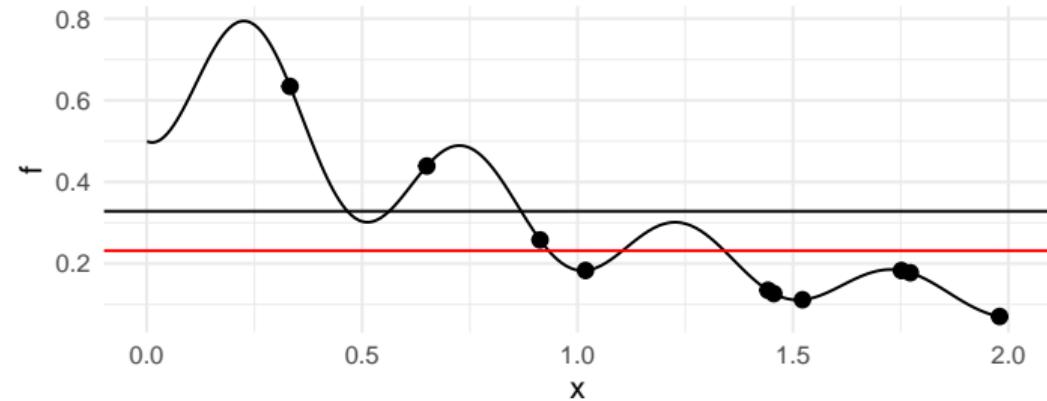
Monte Carlo Importance, $N = 1$



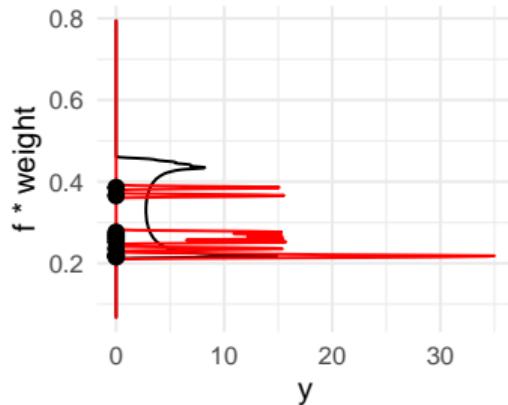
Monte Carlo



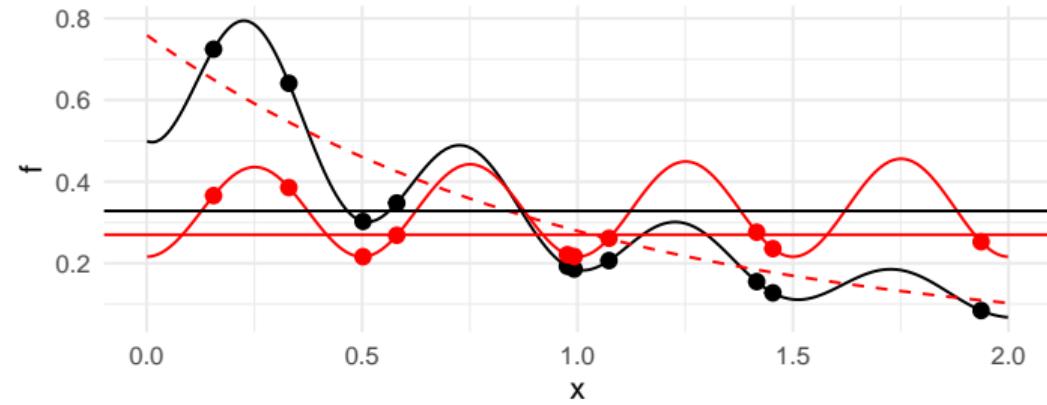
Monte Carlo, $N = 10$



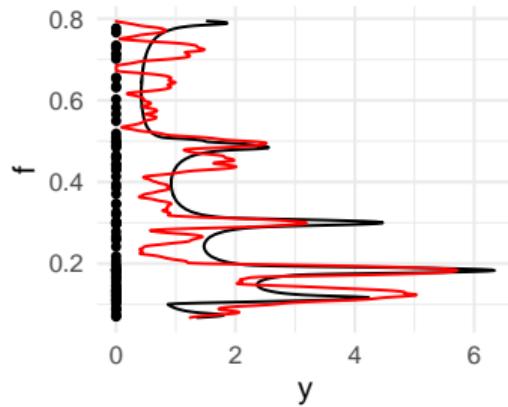
MC Importance



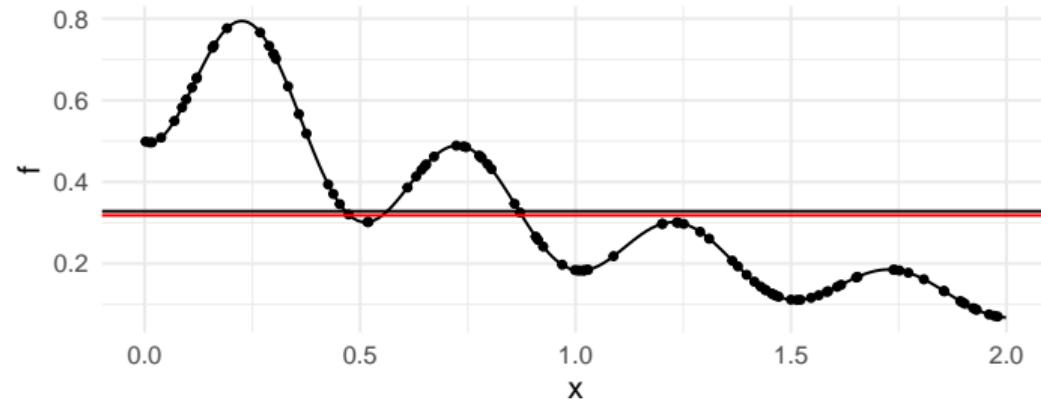
Monte Carlo Importance, $N = 10$



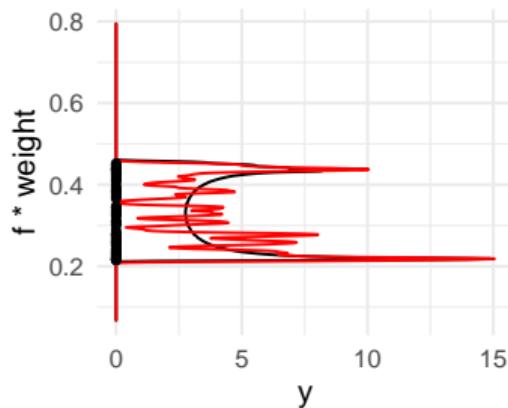
Monte Carlo



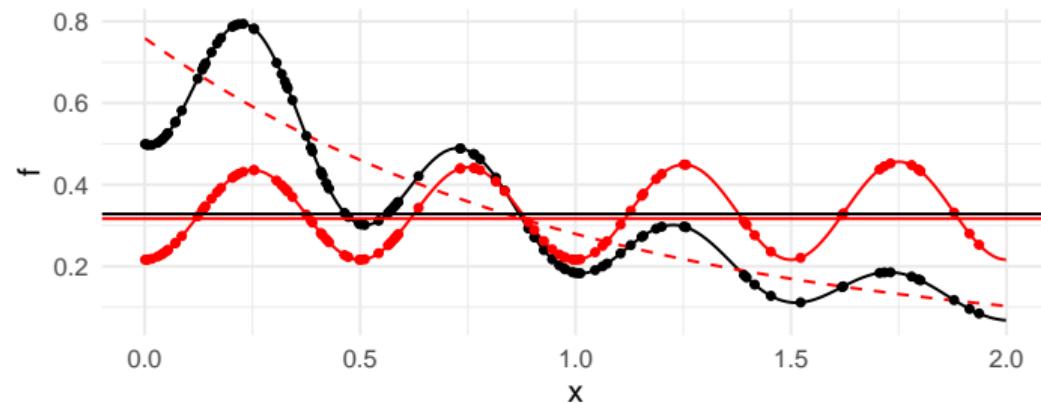
Monte Carlo, $N = 100$



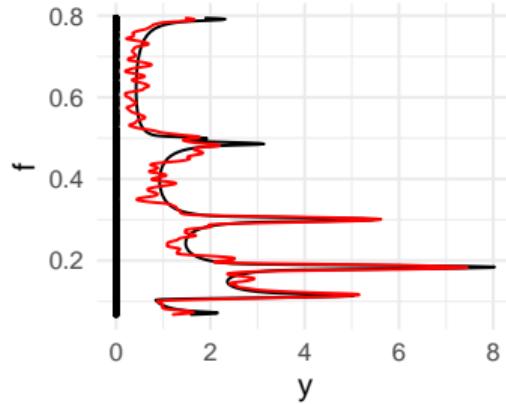
MC Importance



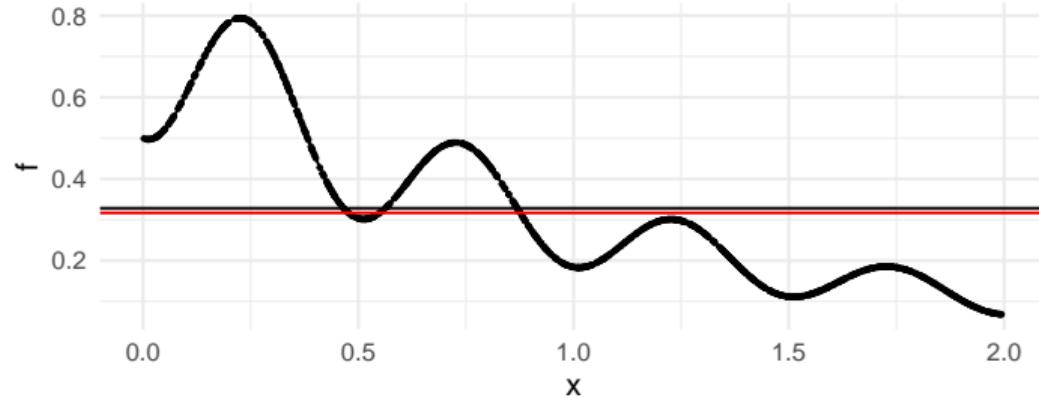
Monte Carlo Importance, $N = 100$



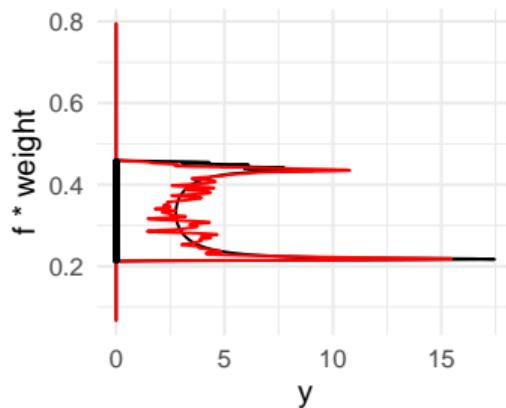
Monte Carlo



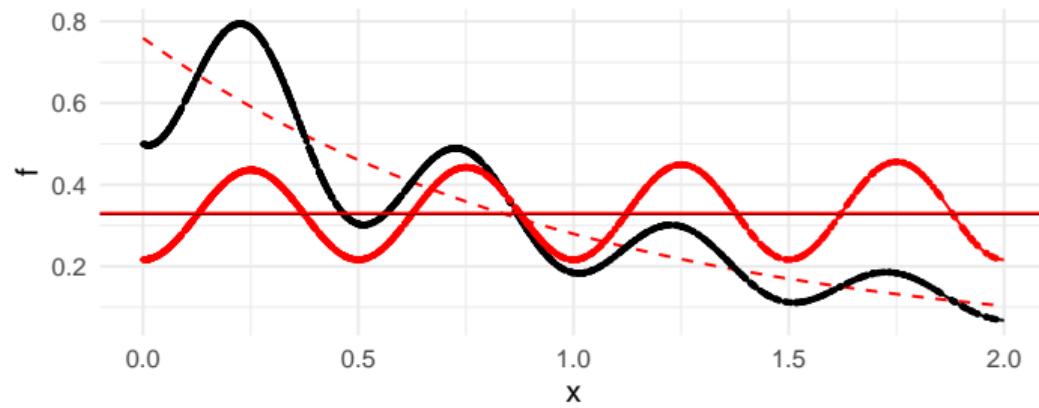
Monte Carlo, $N = 1000$



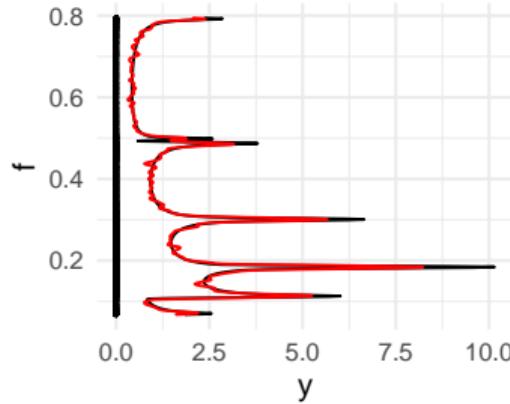
MC Importance



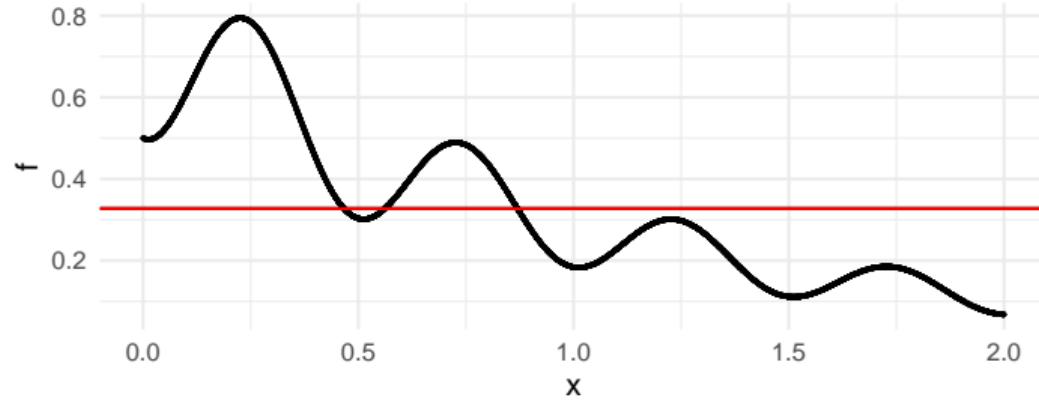
Monte Carlo Importance, $N = 1000$



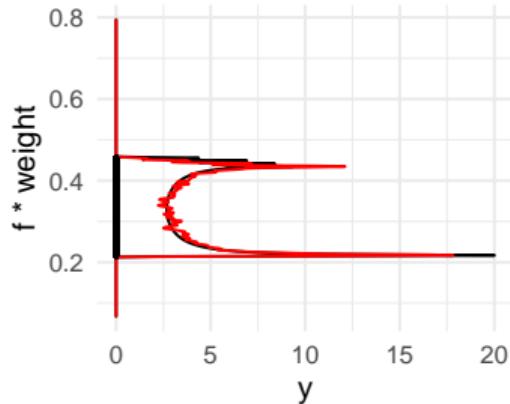
Monte Carlo



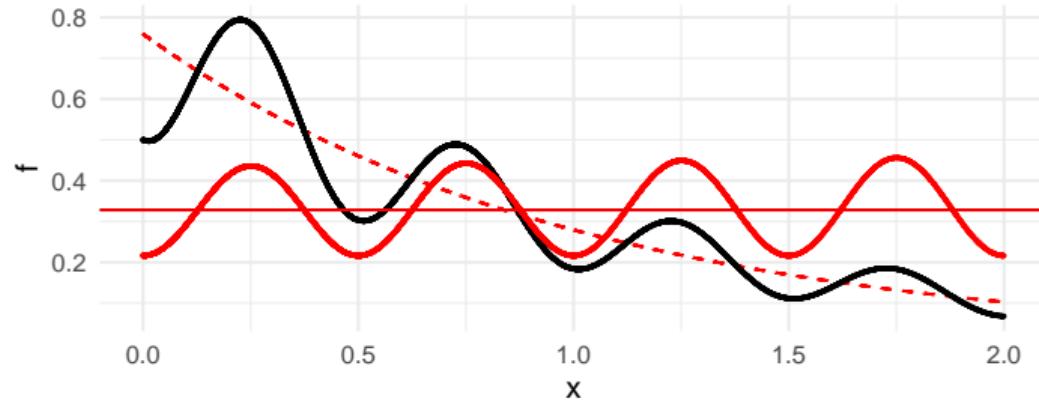
Monte Carlo, $N = 10000$



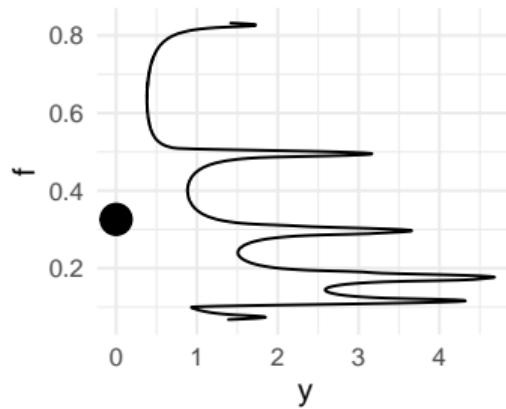
MC Importance



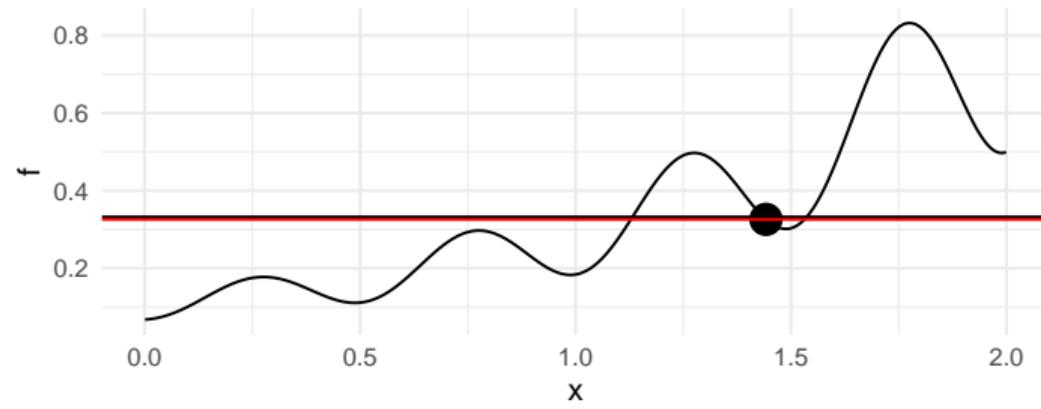
Monte Carlo Importance, $N = 10000$



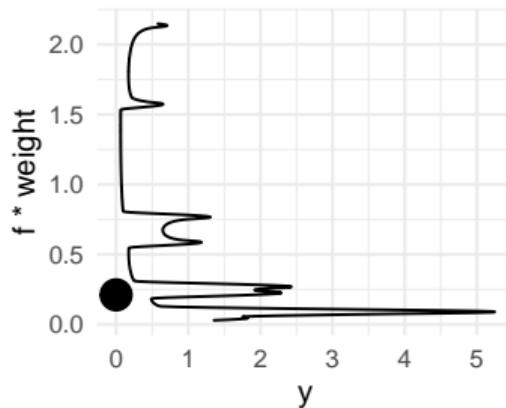
Monte Carlo



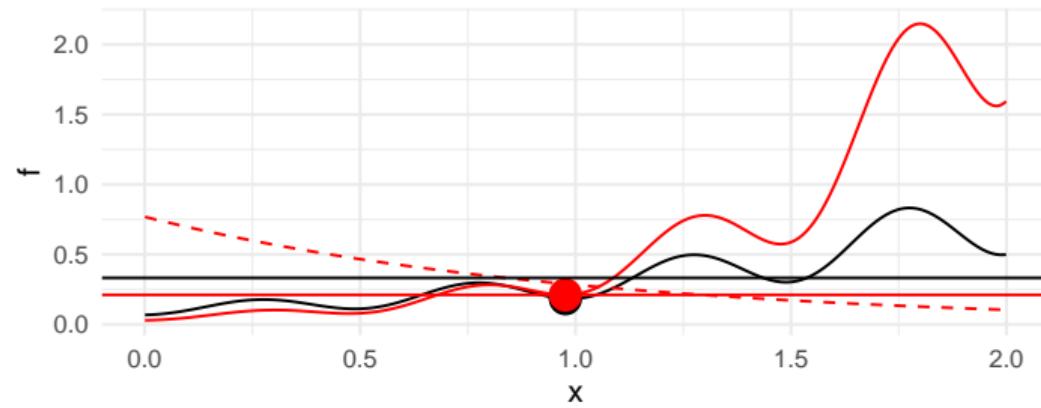
Monte Carlo, $N = 1$



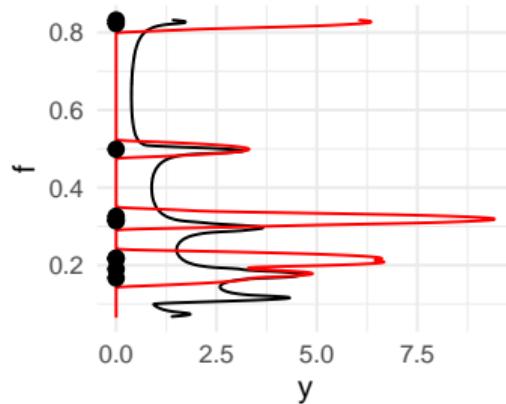
MC Importance



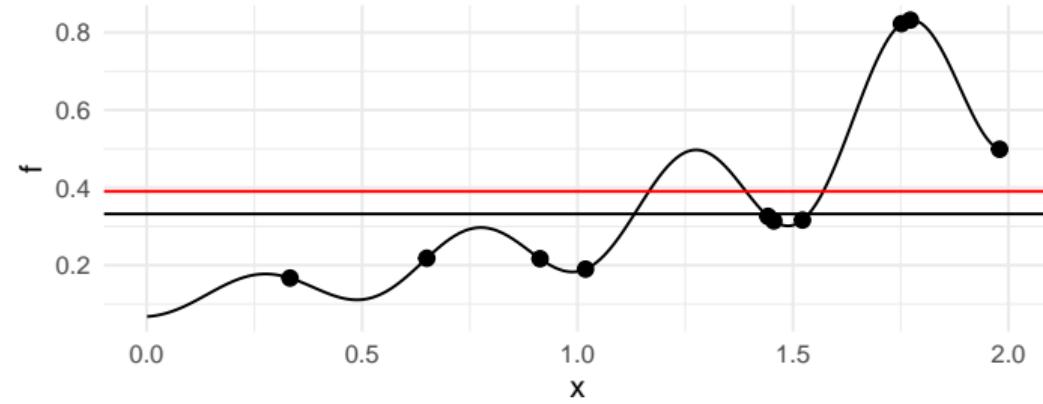
Monte Carlo Importance, $N = 1$



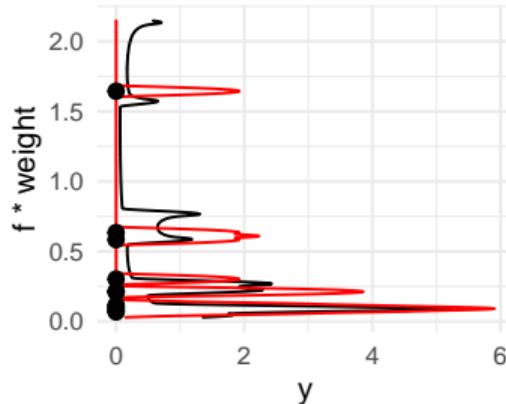
Monte Carlo



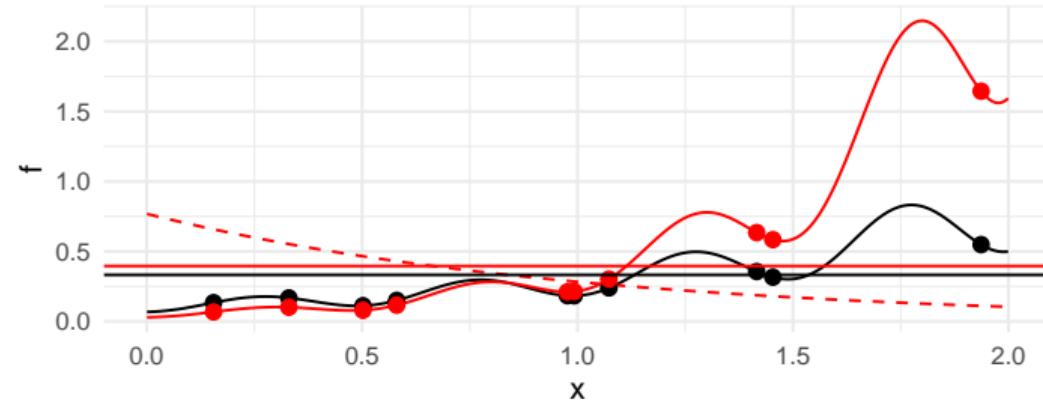
Monte Carlo, $N = 10$



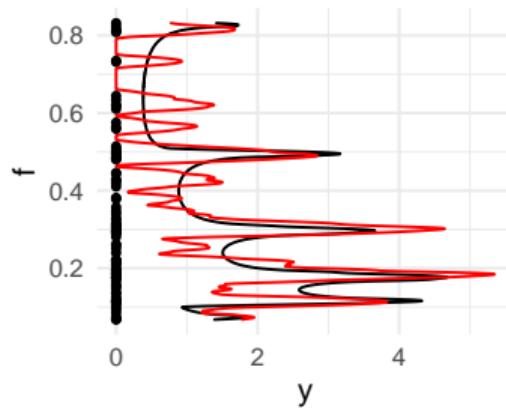
MC Importance



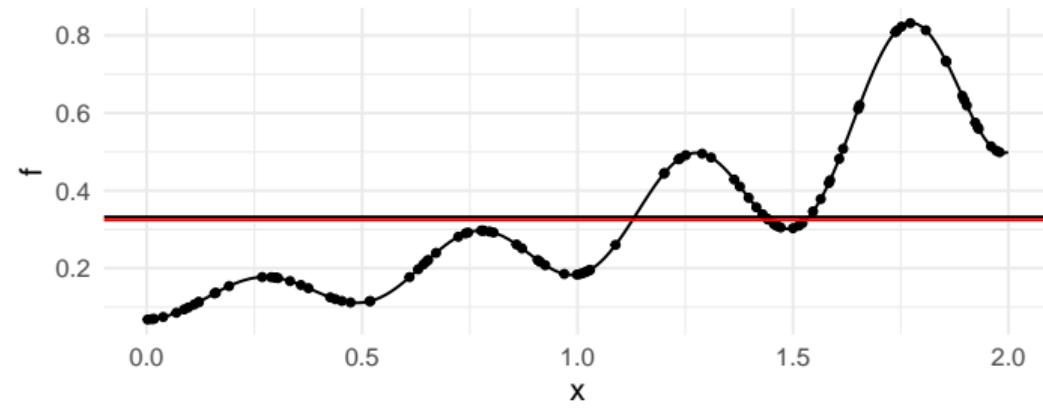
Monte Carlo Importance, $N = 10$



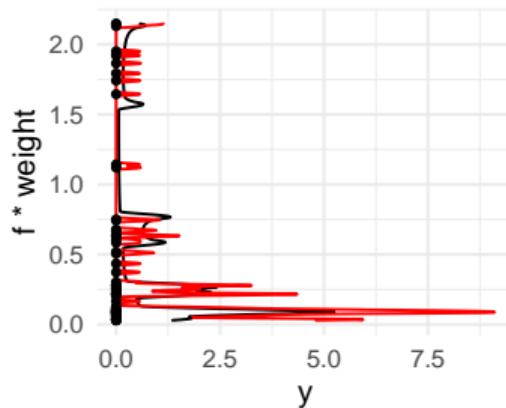
Monte Carlo



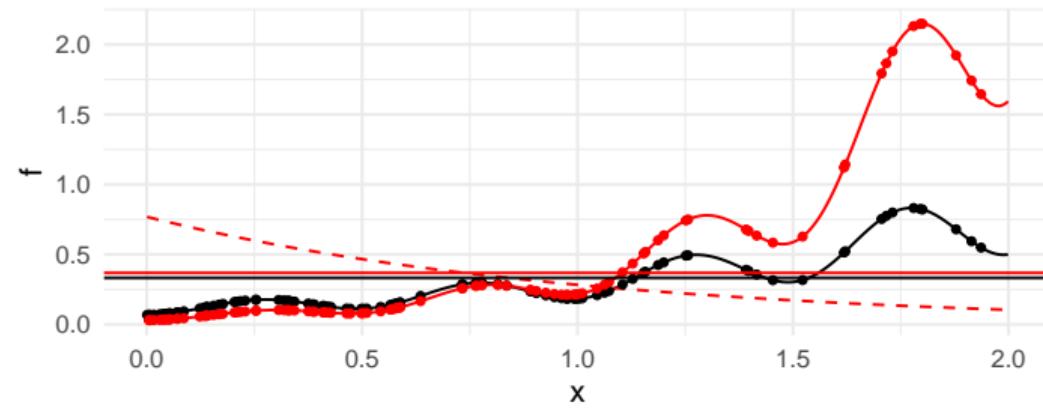
Monte Carlo, $N = 100$



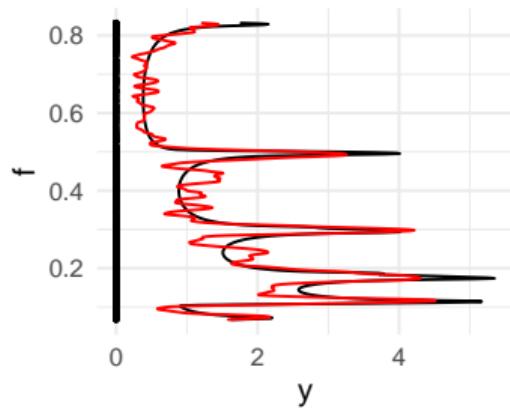
MC Importance



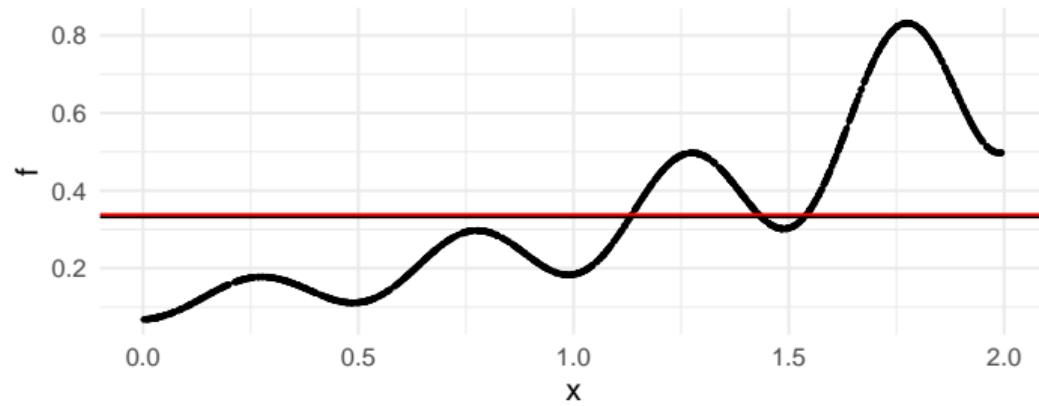
Monte Carlo Importance, $N = 100$



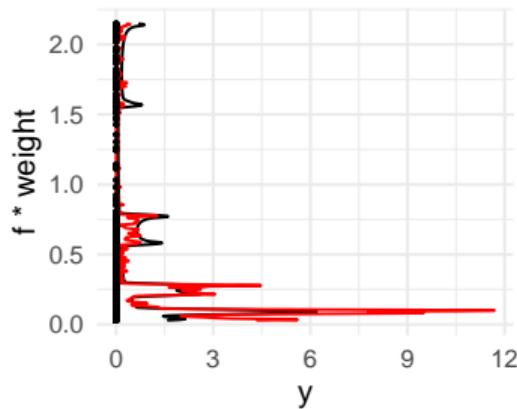
Monte Carlo



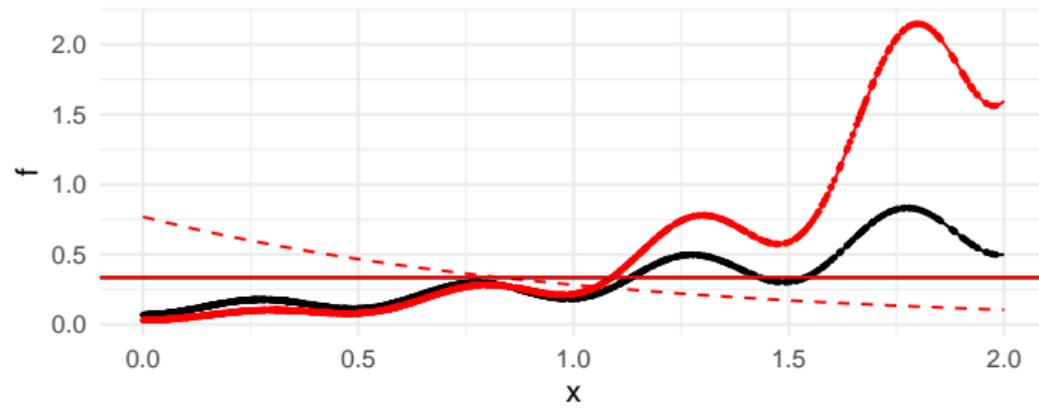
Monte Carlo, $N = 1000$



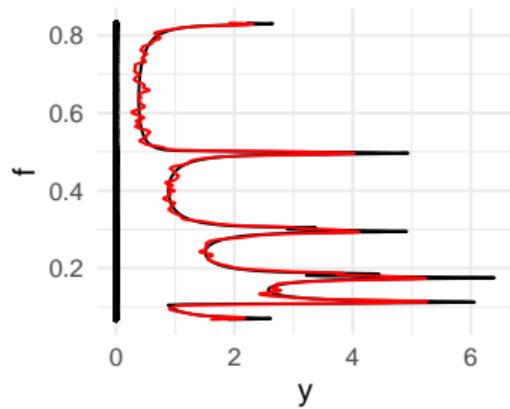
MC Importance



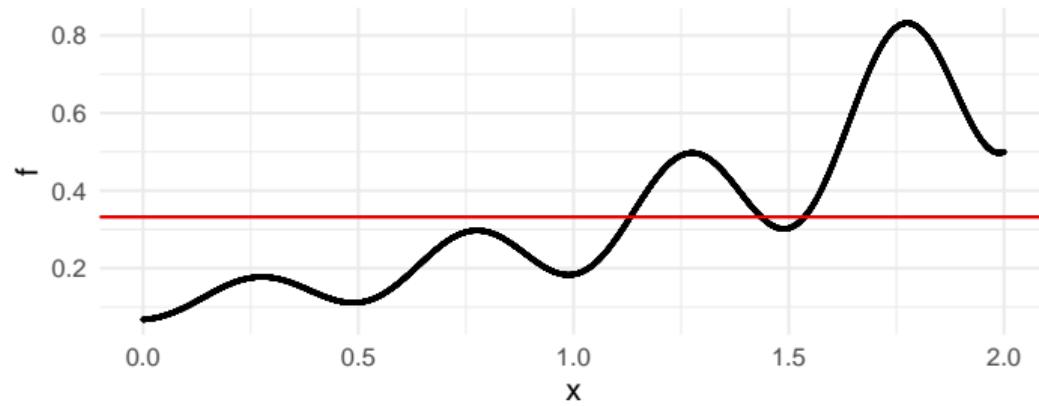
Monte Carlo Importance, $N = 1000$



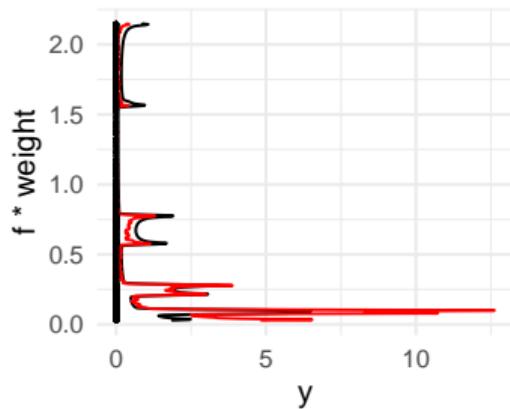
Monte Carlo



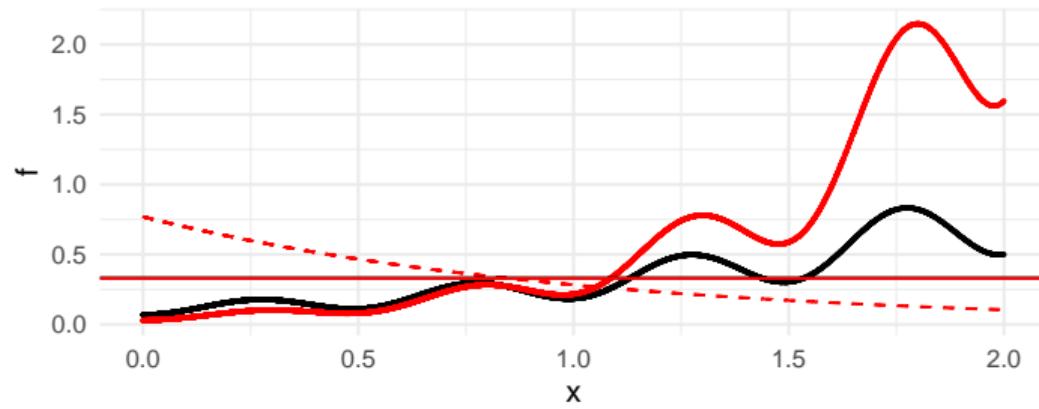
Monte Carlo, $N = 10000$



MC Importance



Monte Carlo Importance, $N = 10000$



Bias/Variance tradeoff

- ▶ The bias for the deterministic Midpoint rule is $\mathcal{O}(h^2)$, where h is the distance between each integration point.
- ▶ In d dimensions, the spacing in regular lattice for a cube $[0, 1]^d$ is $h = N^{-1/d}$, which means that the bias is $\mathcal{O}(N^{-2/d})$.
- ▶ The standard deviation for plain Monte Carlo integration is $\mathcal{O}(N^{-1/2})$.
- ▶ For $1 \leq d < 4$, the bias for the Midpoint rule decreases faster in N than the MC std.dev.
- ▶ For $d > 4$, the MC std.dev. decreases faster.
- ▶ Conclusion: *Monte Carlo integration is relatively efficient for high dimensional integrals.*

Bias/Variance tradeoff

Can we get the best of both worlds? Partly, via *stratified* Monte Carlo:

- ▶ Split the integration domain into N disjoint subdomains: $\Omega = \cup_{k=1}^N \Omega_k$
- ▶ Sample each $\mathbf{x}_k \sim \text{Unif}(\Omega_k)$ and let $\hat{I} = \sum_{k=1}^N |\Omega_k| f(\mathbf{x}_k)$, which is unbiased
- ▶ Take a regular lattice with $|\Omega_k| = |\Omega|/N$, $h = \mathcal{O}(N^{-1/d})$. If f has bounded gradients, then $\text{Var}_{\mathbf{x} \sim \text{Unif}(\Omega_k)} [f(\mathbf{x})] = \mathcal{O}(h^2)$, and

$$\text{Var}_{\{\mathbf{x}_k \sim \text{Unif}(\Omega_k)\}} \left[\sum_{k=1}^N |\Omega_k| f(\mathbf{x}_k) \right] = \dots = \frac{|\Omega|^2}{N^2} N \mathcal{O}(h^2) = \mathcal{O}(N^{-2/d-1})$$

so the std.dev. $\mathcal{O}(N^{-1/d-1/2})$ decays faster in N than for plain MC, and also faster than plain Midpoint integration for $d > 2$. The gain over plain MC diminishes with increasing d .

- ▶ *Quasi Monte Carlo* integration uses a deterministic space filling sequence, converting variance to bias.