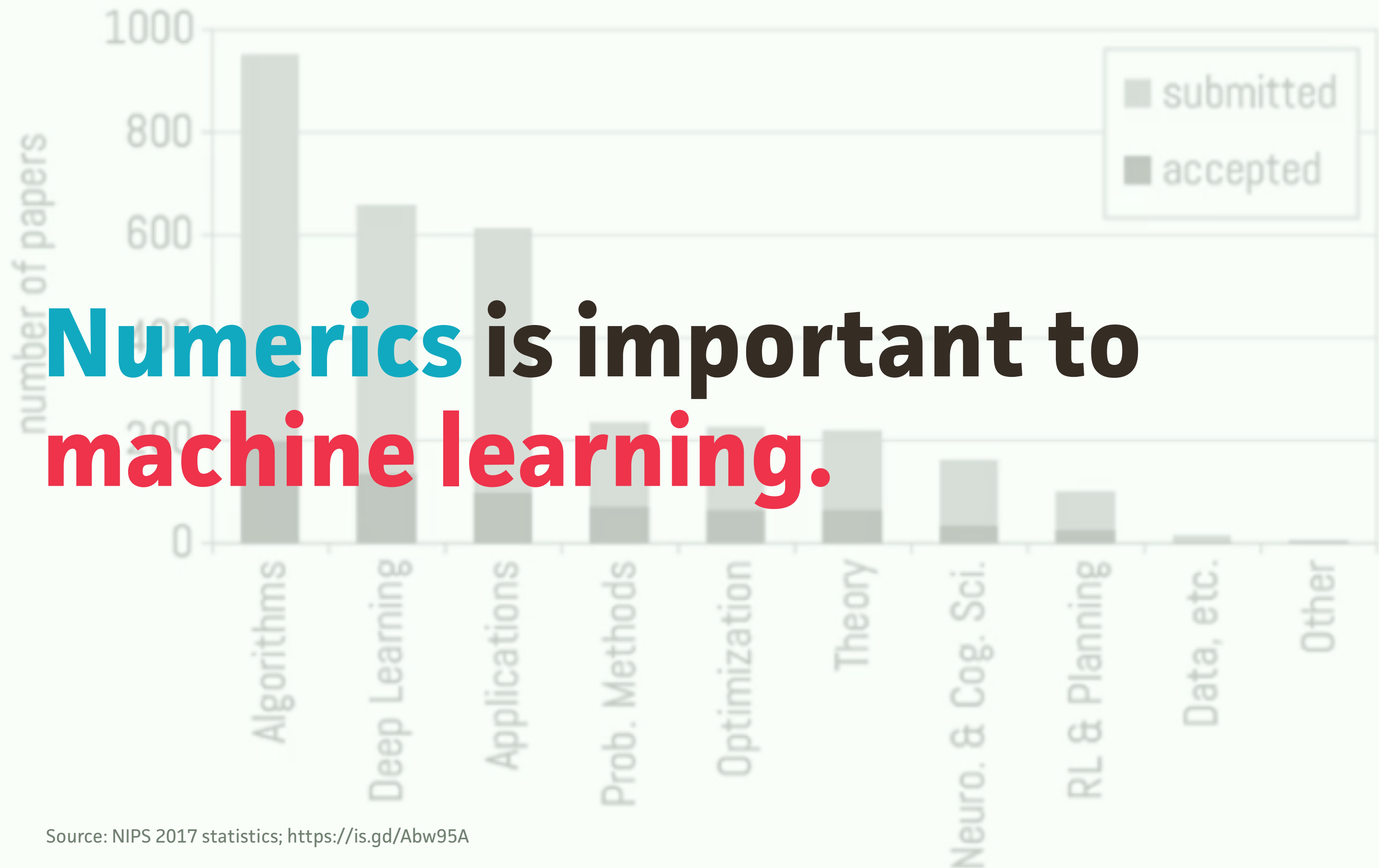# PROBABILISTIC NUMERICS:

## NANO-MACHINE-LEARNING

Michael A Osborne, @maosbot

**Numerics is important to machine learning.**

# Which **numerics problems** have you needed solved in the **last month?**

1. Linear algebra.

2. Optimisation.

3. Global optimisation.

4. Integration.

5. Ordinary differential equations.

C. Nassou
Crecucz Bosch
Wilhelm Inf.
Bernsort
Schwark eck

Admiraliteijt
Inf.

fnz Orange der bestende
Eijs eck
C Diffrigher
C Denis
C Hemskirck
Tern Eck
Troost Eck

ZEMBLA

NOVA

C Plantij.
Loms Baij.
Gros Baij.

Langenes

Ersten Eck.
Schank Eck
Schwark klipp.

Costini sarch.

Creuk Eck.
Schank Eck.
S. laurens Baij.
Meel hauen.
Niduage Inf.
a Insel
Lodigen Eck
Traijebay
Abgoetter Eck

Zwischen Eck
Waigats

TARTARIAE
PARS

Stenden fret.

Oby fluuius

fretum Nassou

De thon Eck

Valcken Eck

Matsse dela

kildun

Olena
cola

Insula

colgoy

Togar

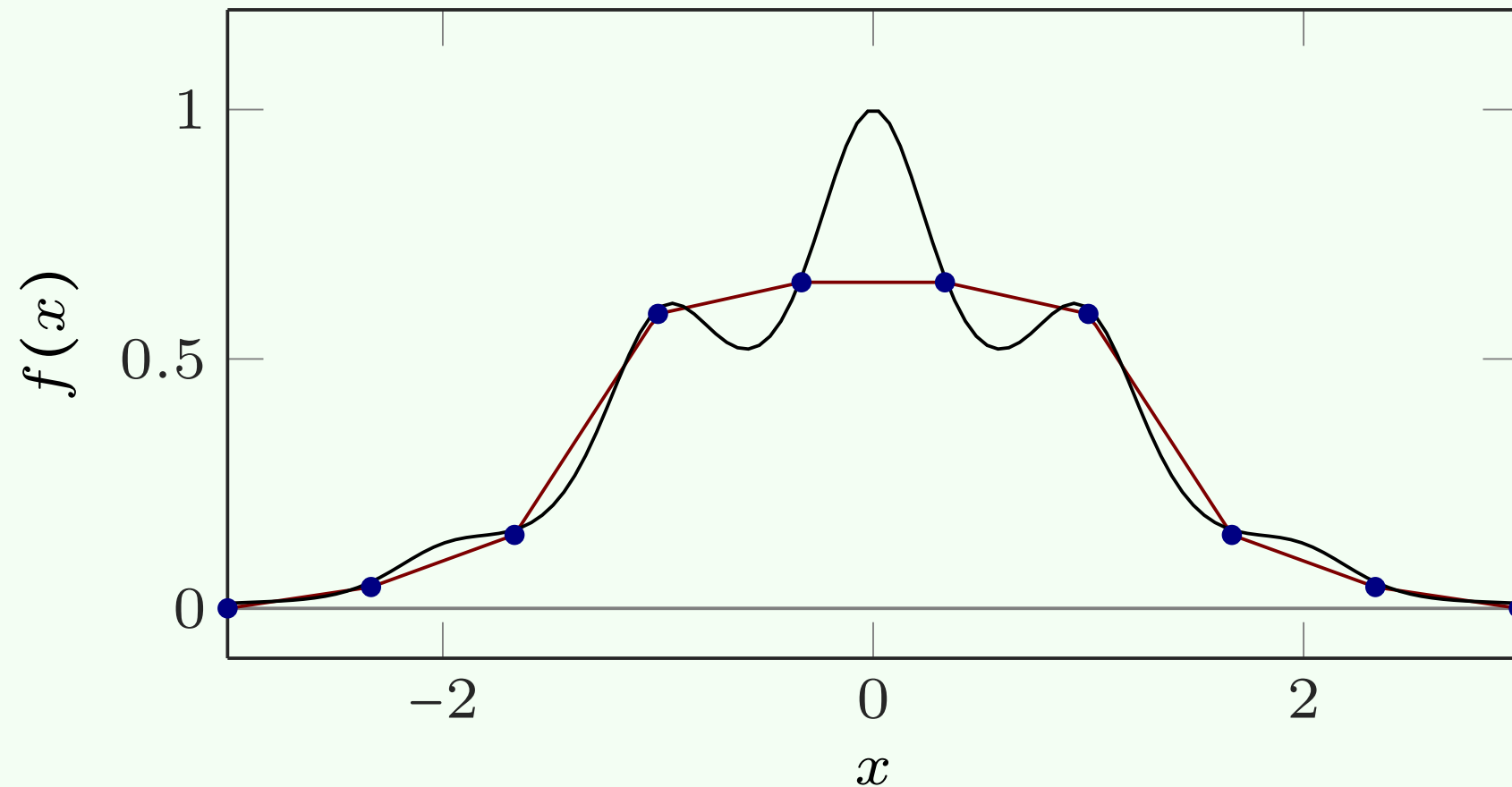# Numerics allows computation-fuelled expeditions beyond the analytic frontier.

# The answer to a numeric problem can only be approximated,

e.g.

$$F = \int_{-3}^{3} f(x)\mathrm{d}x$$

for

$$f(x) = \exp\left(-\left(\sin(3x)\right)^2 - x^2\right).$$

# Machine learning treats algorithms as agents.

# Probabilistic numerics treats numeric algorithms as agents.

```python
import numpy as np
import platform
import subprocess
import nlopt
from sklearn.utils import check_random_state
from scipy.stats import beta, norm



class R        :
    def __init__(self):
        self.name = 'Robotic Arm Simulator'

    def                (        system         )

    def abs_pos(self, jt_angle):
        assert jt_angle.ndim == 1, 'jt_angle has to be one dimensional'
        assert len(jt_angle) == 3, '              has to have 3 inputs'

    if                                  :
        args = str('./robot_arm ' + str(jt_angle[0]) +' '+ str(jt_angle
        proc = subprocess.Popen(args, shell=True, stdout=subprocess.P
    if                windows          :
        args = str('robot_arm.exe ' + str(jt_angle[0]) +' '+ str(jt_angl
        proc = subprocess.Popen(args, stdout=subprocess.PIPE)

        output = proc.stdout
        for line in output:
            output = line
        proc.kill()
        return np.array([float(out) for out in output.split()])
```

As motivation:

1. numeric **error** is significant;

2. numeric methods are **generic**;
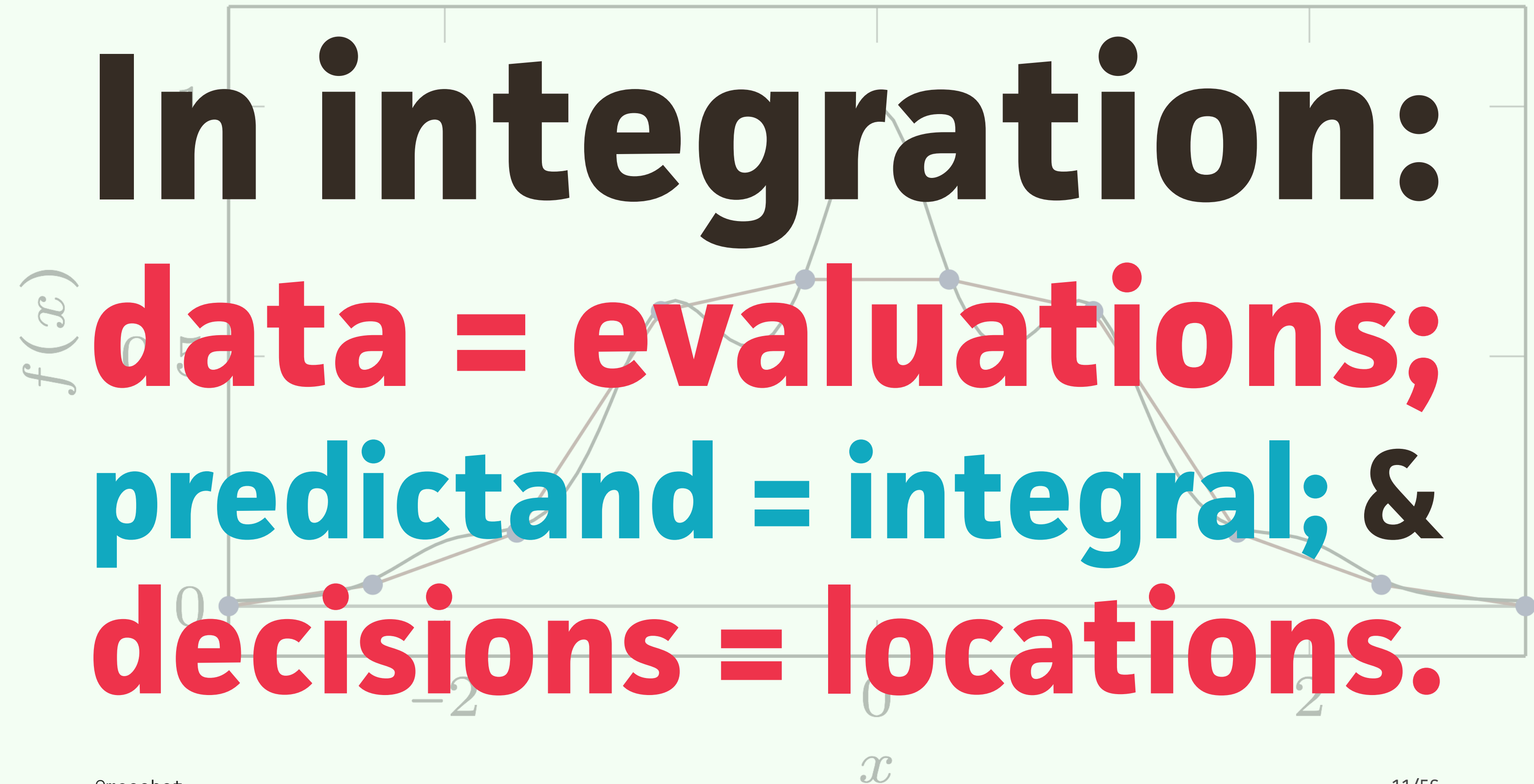
3. our numerics problems **tax our computation.**

# An agent receives data, predicts, & then makes decisions.

# In integration:
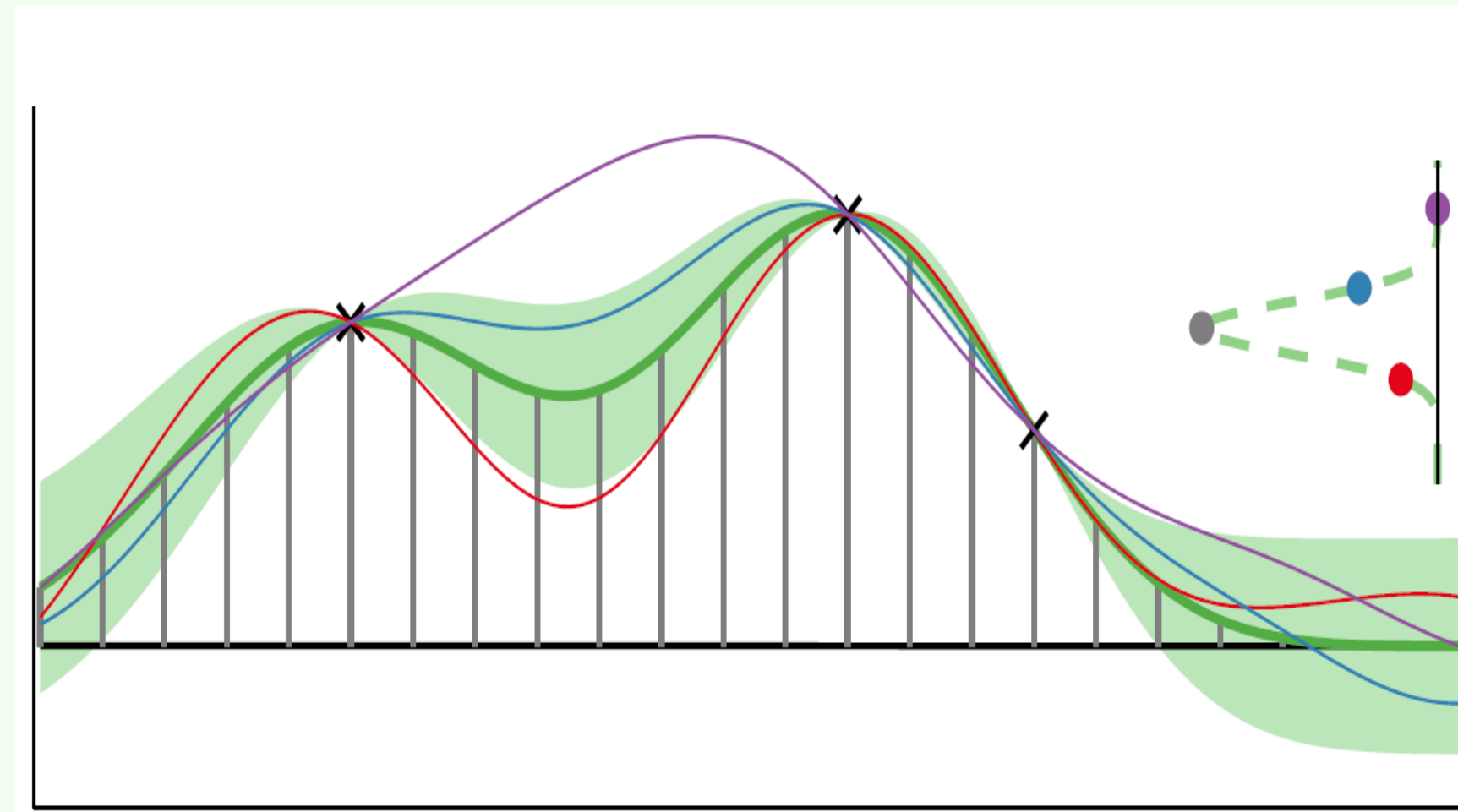## data = ?;
## predictand = ?; &
## decisions = ?.

# In integration:

**data = evaluations;**

**predictand = integral; &**

**decisions = locations.**

# Bayesian quadrature is probabilistic numerics for integration.

# An agent is defined by its **prior** and loss function.

# With a **Gaussian process** prior for the integrand, the **integral is joint Gaussian.**

The trapezoidal rule is the posterior mean estimate for the integral

$$F = \int_a^b f(x)\,\mathrm{d}x$$

under any centered Wiener process prior

$$p(f) = \mathcal{GP}(f; 0, k)$$
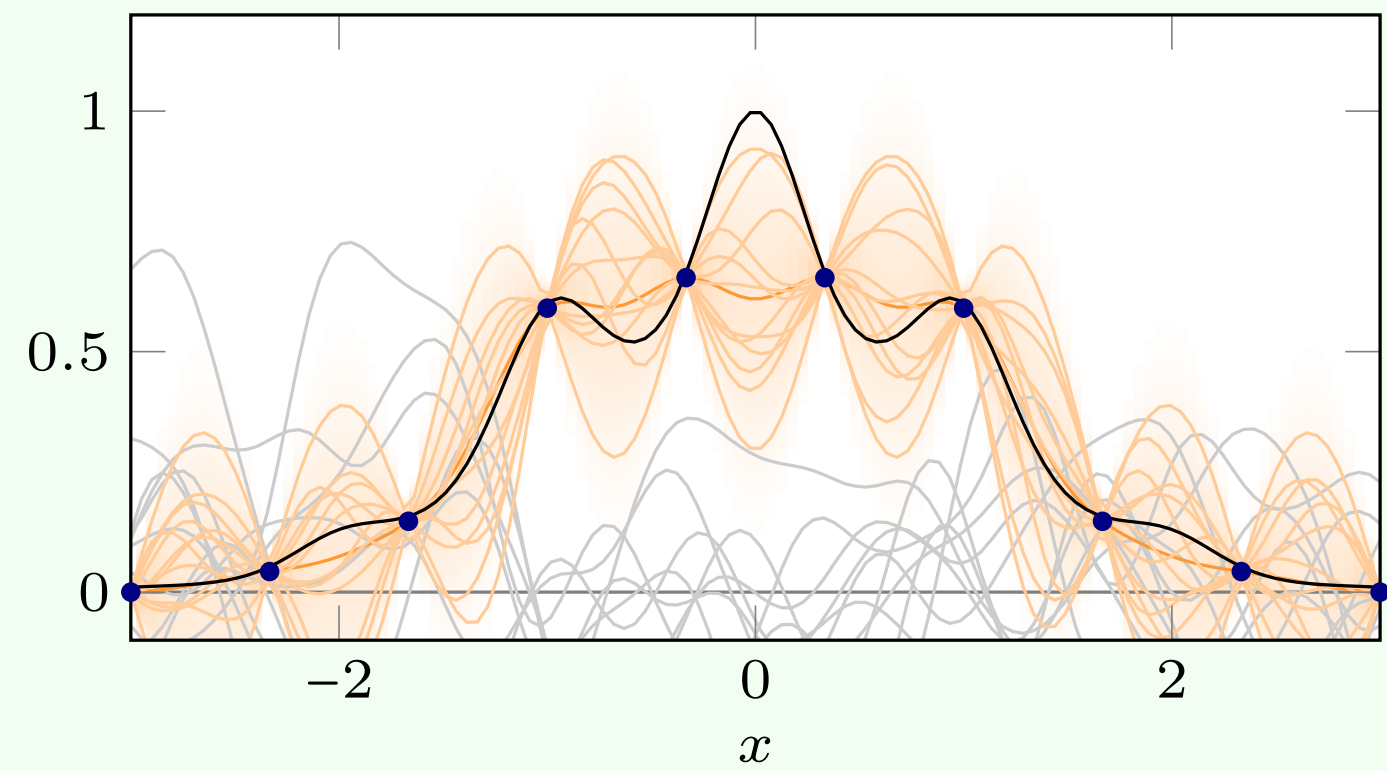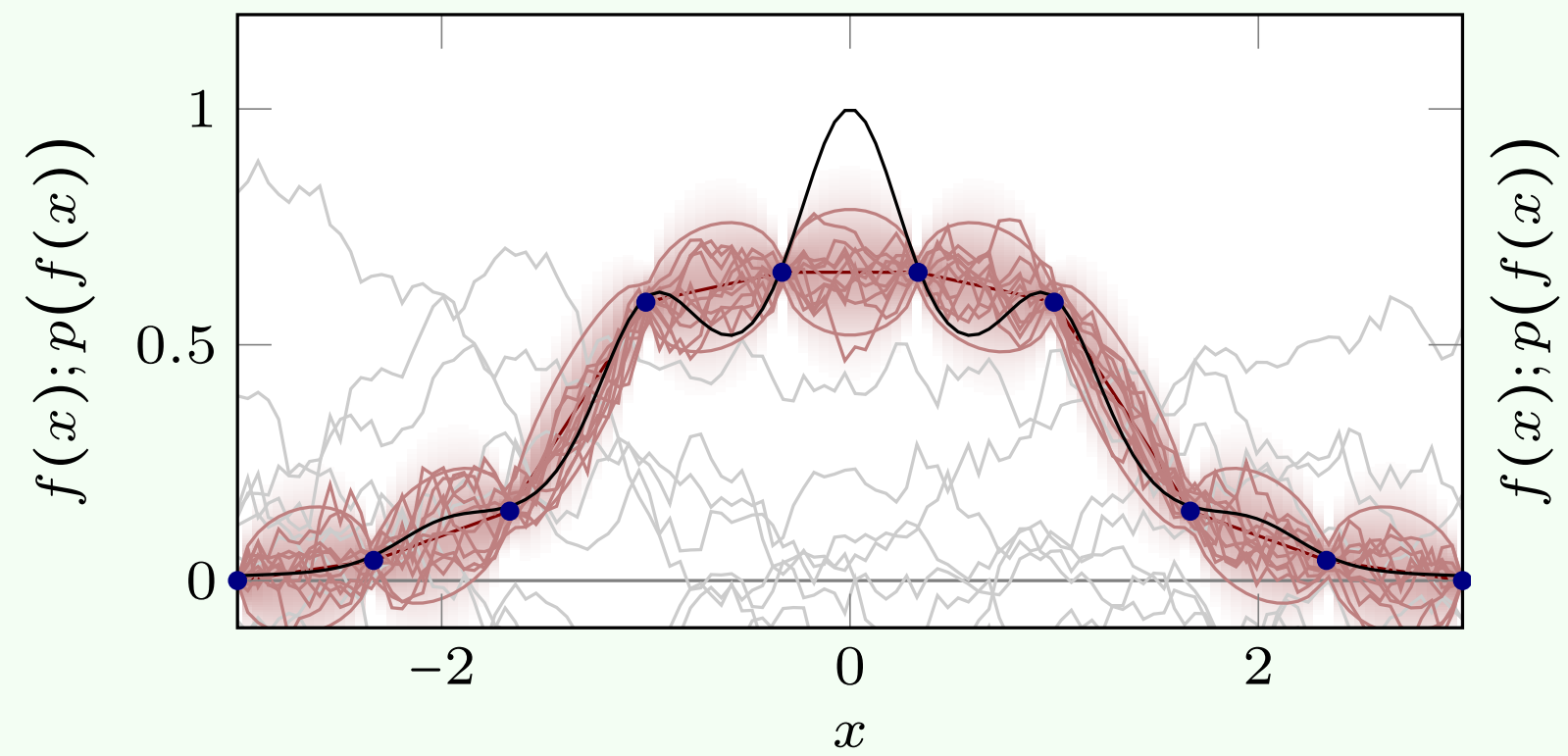
with

$$k(x, x') = \theta^2 \left(\min(x, x') - \chi\right)$$

for arbitrary $\theta \in \mathfrak{R}_+$ and $\chi < a \in \mathfrak{R}$.

1 **procedure** $\text{INTEGRATE}(@f, a, b, N, \theta)$

2 $\quad \delta := (b - a)/(N - 1)$      // choose step size

3 $\quad x \leftarrow a, y_1 = f(a), m \leftarrow 0, v \leftarrow 0,$      // initialise

4 $\quad$ **for** $i = 2, \ldots, N$ **do**

5 $\quad\quad x \leftarrow x + \delta$      // step

6 $\quad\quad y_i \leftarrow f(x)$      // evaluate

7 $\quad\quad m \leftarrow m + \delta/2(y_{i-1} + y_i)$      // update estimate

8 $\quad\quad v \leftarrow v + \delta^3/12$      // update error estimate

9 $\quad$ **end for**

10 $\quad$ **return** $\mathbb{E}(F) = m, \text{var}(F) = \theta^2 v$      // probabilistic output

11 **end procedure**

**procedure** INTEGRATE($@f, a, b, N, \theta$)
1

$\delta := (b - a)/(N - 1)$      // choose step size
2

$x \leftarrow a, y_1 = f(a), m \leftarrow 0, v \leftarrow 0,$      // initialise
3

**for** $i = 2, \ldots, N$ **do**
4

$x \leftarrow x + \delta$      // step
5

$y_i \leftarrow f(x)$      // evaluate
6

$m \leftarrow m + \frac{\delta}{2}(y_{i-1} + y_i)$      // update estimate
7

$v \leftarrow v + \delta^3/12$      // update error estimate
8

**end for**
9

**return** $\mathbb{E}(F) = m, \text{var}(F) = \theta^2 v$      // probabilistic output
10

**end procedure**
11

# The trapezoid rule **is** Bayesian quadrature.

# Quiz: The convergence rate of the trapezoid rule is $\mathcal{O}(N^{-1})$: what is the rate of Monte Carlo?

1.    $\mathcal{O}\big(\exp(-N)\big)$

2.    $\mathcal{O}\Big(\exp\big(-N^{-\frac{1}{2}}\big)\Big)$

3.    $\mathcal{O}\big(N^{-1}\big)$

4.    $\mathcal{O}\big(N^{-\frac{1}{2}}\big)$

# Quiz: The convergence rate of the trapezoid rule is $\mathcal{O}(N^{-1})$: what is the rate of Monte Carlo?

1. $\quad \mathcal{O}\big(\exp(-N)\big)$

2. $\quad \mathcal{O}\big(\exp\big(-N^{-\frac{1}{2}}\big)\big)$

3. $\quad \mathcal{O}(N^{-1})$

4. $\quad \mathcal{O}\big(N^{-\frac{1}{2}}\big)$ – **arguably the worst possible rate.**

# Monte Carlo **is also** Bayesian quadrature.

The Monte Carlo estimate

$$\int f(x)\, p(x)\, \mathrm{d}x \simeq \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$

is maximum a-posteriori under the (improper) prior

$$p(f) = \lim_{c \to 0} \mathcal{GP}\big(0, \theta^2 \mathbb{I}(x = x') + c^{-1}\big)$$

for $\mathbb{I}$ the indicator function and with arbitrary $\theta \in \mathfrak{R}_+$. The corresponding posterior standard deviation estimate on the integral, $\theta\,(b - a)\big/\sqrt{N}$, matches the convergence rate of the Monte Carlo estimator.

Monte Carlo

Trapezoidal

Gauss-Legendre

$|F - \hat{F}|$

# samples

**Quadrature is often required to manage model parameters.**

# Managing parameters $\theta$ requires the model evidence,

$$p(\text{data}) = \int p(\text{data} \mid \theta)\, p(\theta)\mathrm{d}\theta.$$

# QUADRATURE
# IS HARD.

Parameter, θ

# Optimisation (maximum likelihood, training) is often used in the place of quadrature.

This approximates as $p(\text{data}) \simeq \int p(\text{data} \mid \theta)\,\delta(\theta - \theta_{\max})\,\mathrm{d}\theta.$

If optimising, **flat optima** are often a better representation of the integral than **narrow optima.**

# **Monte Carlo** has revolutionised Bayesian inference.

Monte Carlo estimators, $\int f(x)\,p(x)\,\mathrm{d}x \simeq \frac{1}{N}\sum_{i=1}^{N} f(x_i),$

**ignore relevant information.**



O'Hagan, A. (1987). Monte Carlo is Fundamentally Unsound. Journal of the Royal Statistical Society. Series D (The Statistician).

# An agent is defined by its prior and loss function.

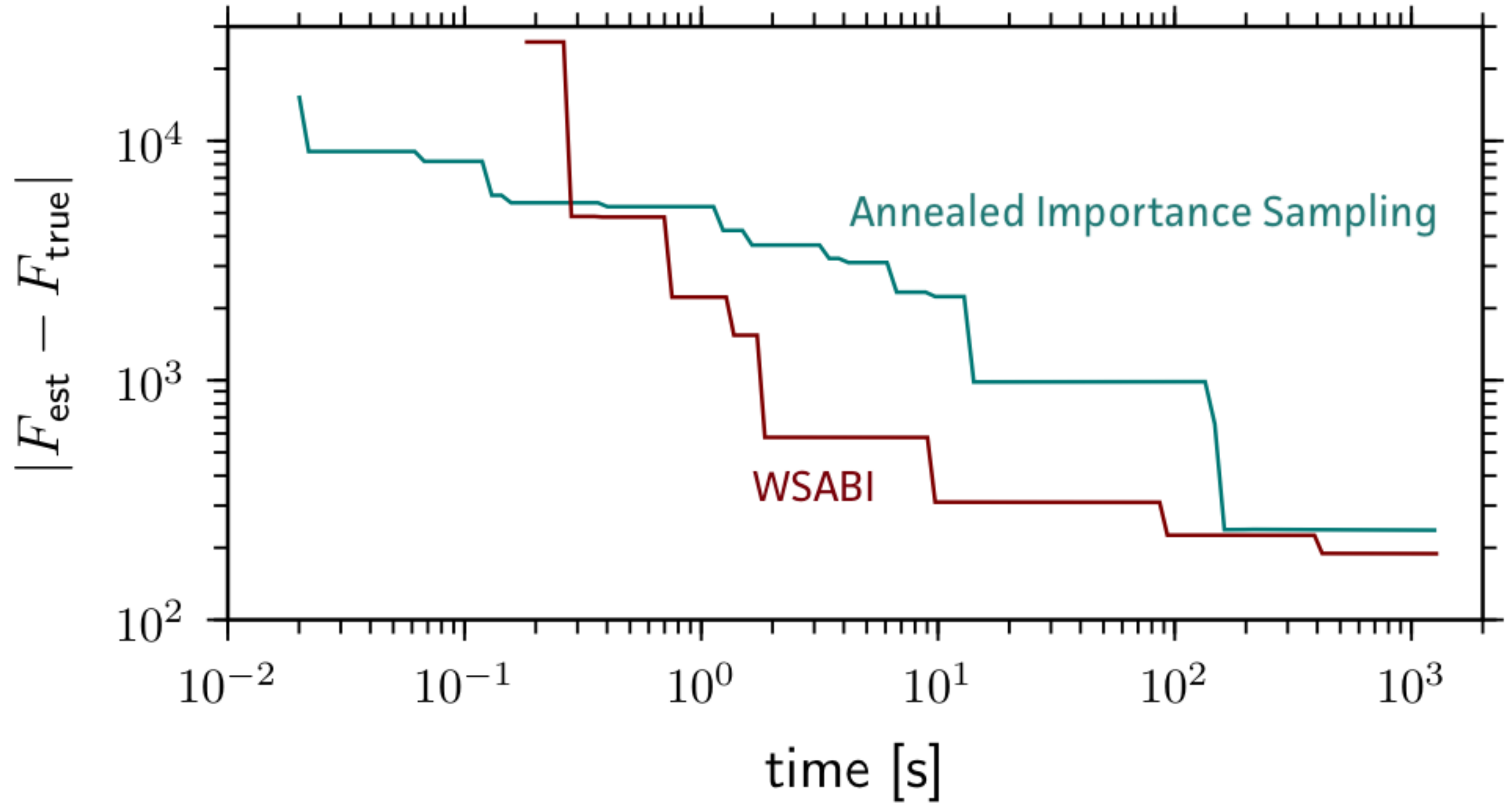A natural **loss function** for quadrature is the **uncertainty in the integral.**

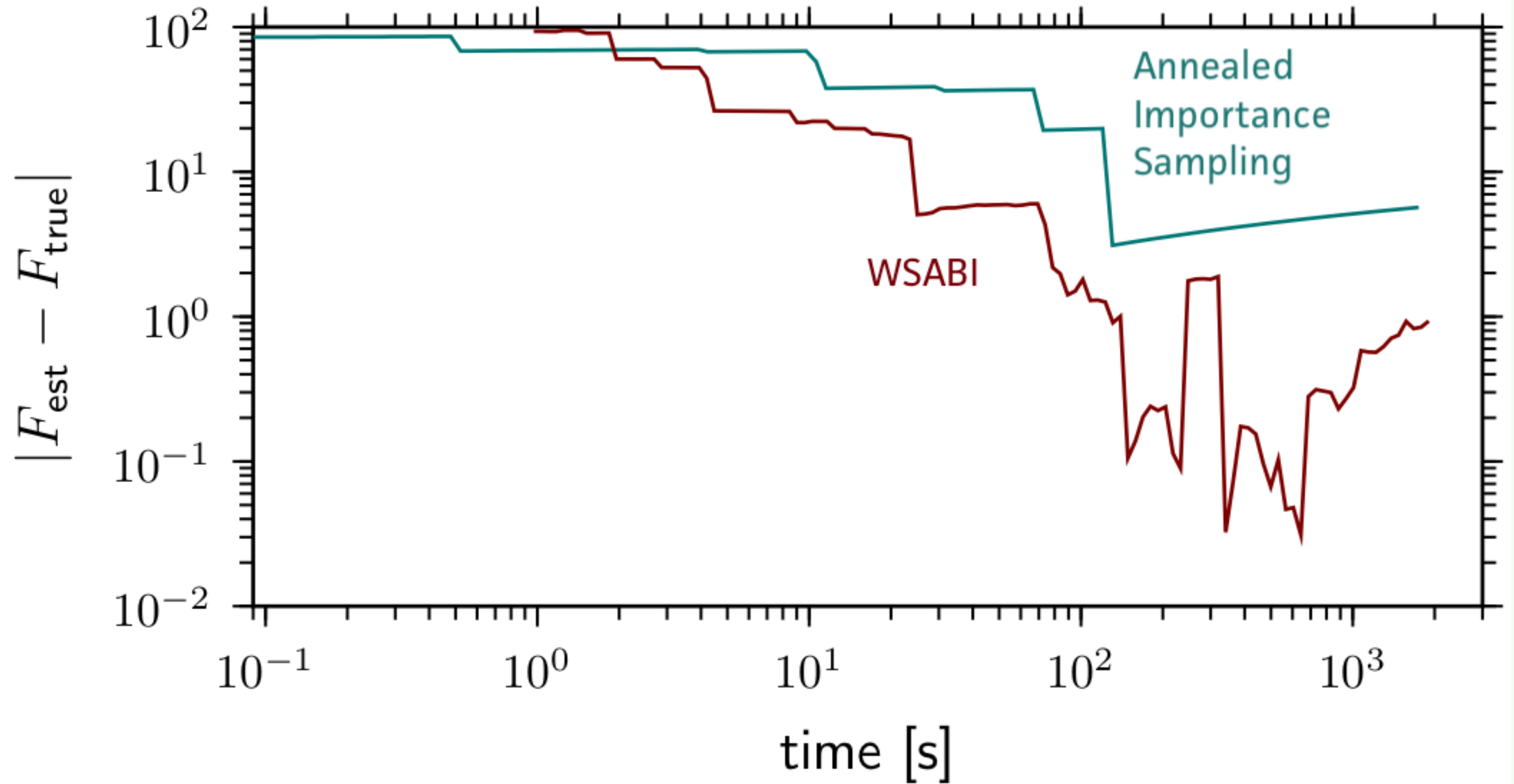**WSABI** uses a loss that is the uncertainty in the integrand.



Source: Gunter, Osborne, Garnett, Hennig, & Roberts (2014). Sampling for Inference in Probabilistic Models with Fast Bayesian Quadrature. NIPS.

synthetic (moG)

# yacht hydrodynamics

GP classification, graph

# Overhead can set you free.

# A RANDOM NUMBER IS A DECISION.

# A random number assumes a completely flat expected loss.

A random number is a **heuristic for a sequence of exploratory decisions:** but why assume **negligible memory?**

# Using random numbers makes your algorithm **unimprovable.**

A random number may be intended to be **unbiased,** but there are **no adversaries in numerics.**

Source: Brown, Mané, Roy, Abadi & Gilmer (2017) "Adversarial Patch".

# Quiz: which of these sequences is random?

1. 6224441111111111444443333333

2. 1693993751058209749445923078

3. 7129042634726105902083360448

4. 1000111111011111111001010000
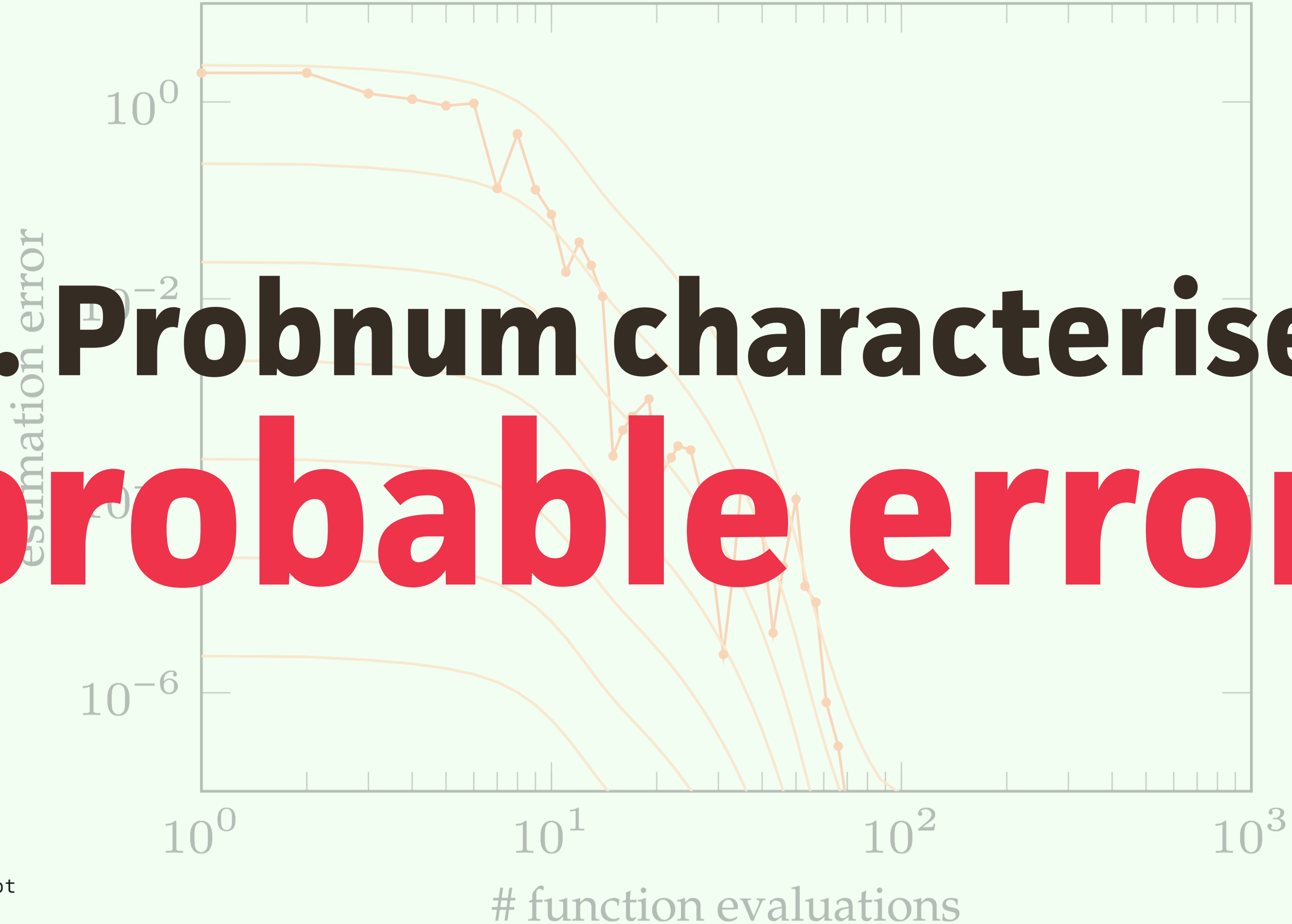
# Quiz: which of these sequences is random?

1. 6224441111111114444443333333: seven d6 rolls with $i$ repeats of the $i$th roll.

2. 1693993751058209749445923078: the 41st to 70th digits of $\pi$.

3. 7129042634726105902083360448: this sequence was generated by the von Neumann method with seed 908344.

4. 1000111111011111111001010000: digits taken from a CD-ROM published by George Marsaglia.

A finite string of random numbers is **encoding some bias!**

# Recall:

## 1. numeric <span style="color:red">error</span> is significant;

## 2. numeric methods are <span style="color:teal">generic</span>;

## 3. our numerics problems <span style="color:red">tax our computation.</span>

# 1. Probnum characterises probable error.

# 2. Probnum tailors procedures to problems.

# 3. Probnum makes better use of computation.

# PROBABILISTIC-NUMERICS.ORG

Numerical algorithms, such as methods for the num
differential equations, as well as optimization algor
They estimate the value of a latent, intractable qua

# probnum.org

Numerical algorithms, such as methods for the num

differential equations, as well as optimization algor

They estimate the value of a latent, intractable qua

# LITERATURE

This page collects literature on all areas of probab
not hesitate to contact us. The fastest way to get
file in /_bibliography, then either send us a pull-re

**QUICK-JUMP LINKS:**

- General and Foundational
- Quadrature
- Linear Algebra
- Optimization
- Ordinary Differential Equations
- Partial Differential Equations

@maosbot

# Huge thanks to Philipp Hennig.