



Machine Learning and the Physical World

Lecture 9 : Probabilistic Numerics

Carl Henrik Ek - che29@cam.ac.uk

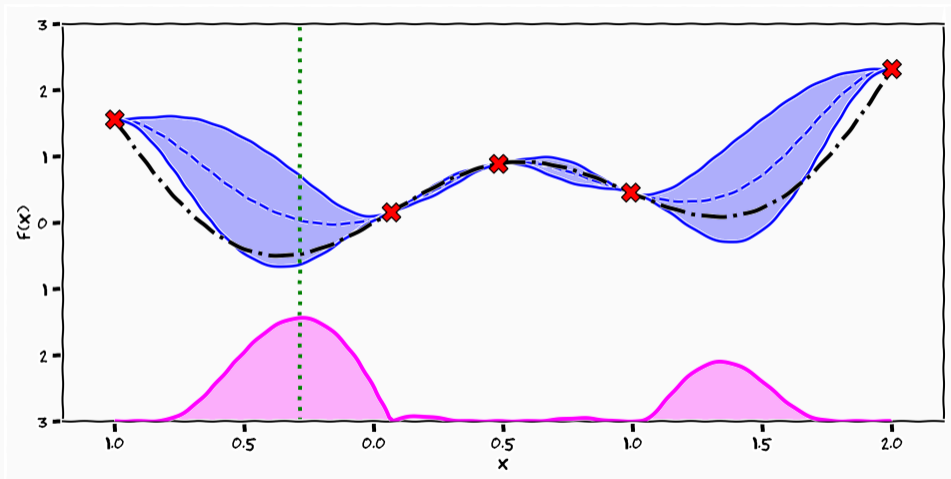
7th of November, 2024

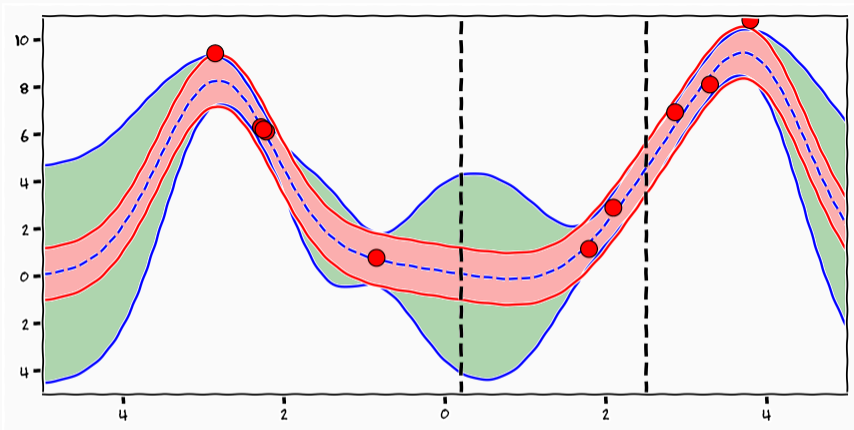
<http://carlhenrik.com>

What do we do with uncertainty?



Bayesian Optimisation





$$y_i = f_i + \epsilon$$

"The need for probability only arises out of uncertainty: It has no place if we are certain that we know all aspects of a problem. But our lack of knowledge also must not be complete, otherwise we would have nothing to evaluate. There is thus a spectrum of degrees of uncertainty. While the probability for the sixth decimal digit of a number in a table of logarithms to equal 6 is 1/10 a priori, in reality, all aspects of the corresponding problem are well determined, and, if we wanted to make the effort, we could find out its exact value. The same holds for interpolation, for the integration methods of Cotes or Gauss, etc"

– Henri Poincare, 1896

Aleatoric/Stochastic "Randomness" inherent in system, or noise in our measurement of system

Aleatoric/Stochastic "Randomness" inherent in system, or noise in our measurement of system

Epistemic Uncertainty related to our ignorance of a the underlying system

Aleatoric/Stochastic "Randomness" inherent in system, or noise in our measurement of system

Epistemic Uncertainty related to our ignorance of a the underlying system

Computational *Uncertainty related to finite computation, or intractable computations*

Data + Model $\xrightarrow{\text{Compute}}$ Prediction

- Computation is expensive, how much knowledge will I gain from computing more?

- Computation is expensive, how much knowledge will I gain from computing more?
- What should I compute in order to reduce my uncertainty as much as possible?

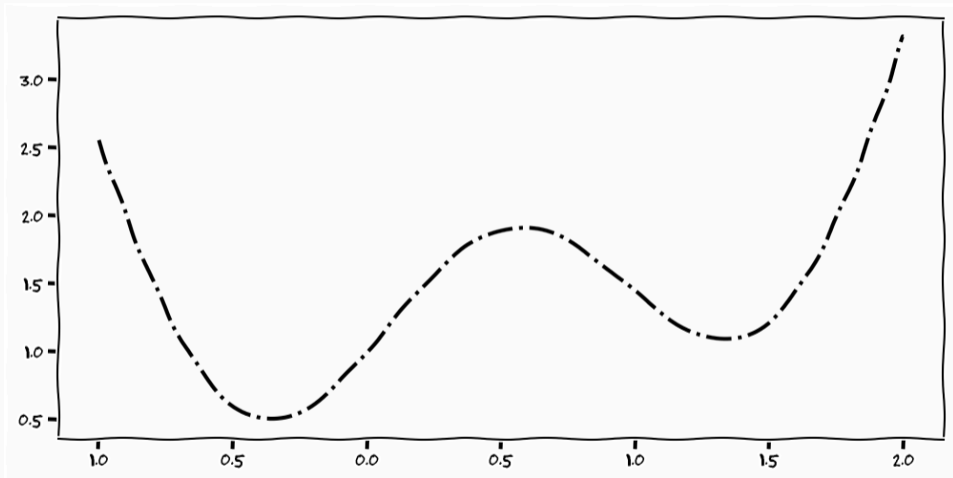
- Computation is expensive, how much knowledge will I gain from computing more?
- What should I compute in order to reduce my uncertainty as much as possible?
- How much should I trust the computation I have done?

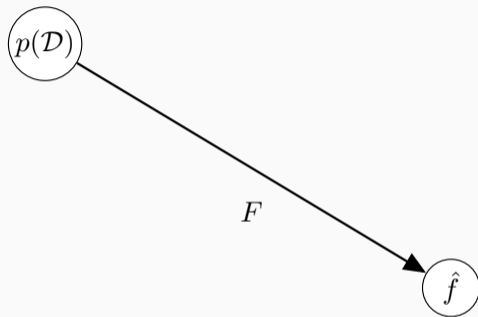
- Computation is expensive, how much knowledge will I gain from computing more?
- What should I compute in order to reduce my uncertainty as much as possible?
- How much should I trust the computation I have done?
- How precise should I do down-stream tasks based on the information from a specific computation?

"[round-off errors] are strictly very complicated but uniquely defined number theoretical functions [of the inputs], yet our ignorance of their true nature is such that we best treat them as random variables."

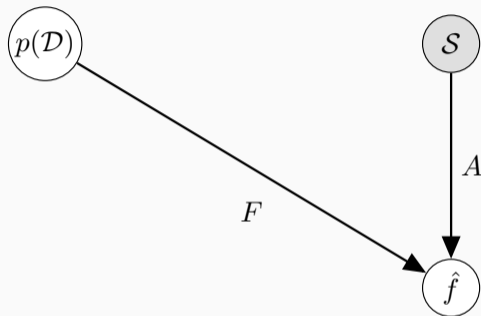
– Neumann et al., [1947](#)

I believe in...

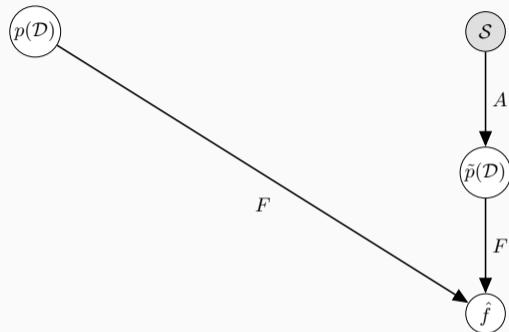




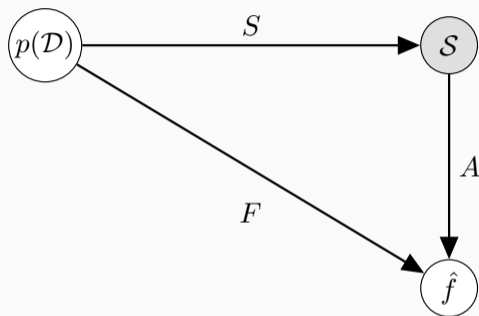
$$F : p(\mathcal{D}) \rightarrow p(\mathcal{Y}|\mathcal{X})$$



$$A \circ \mathcal{S} \approx F \circ p(\mathcal{D})$$

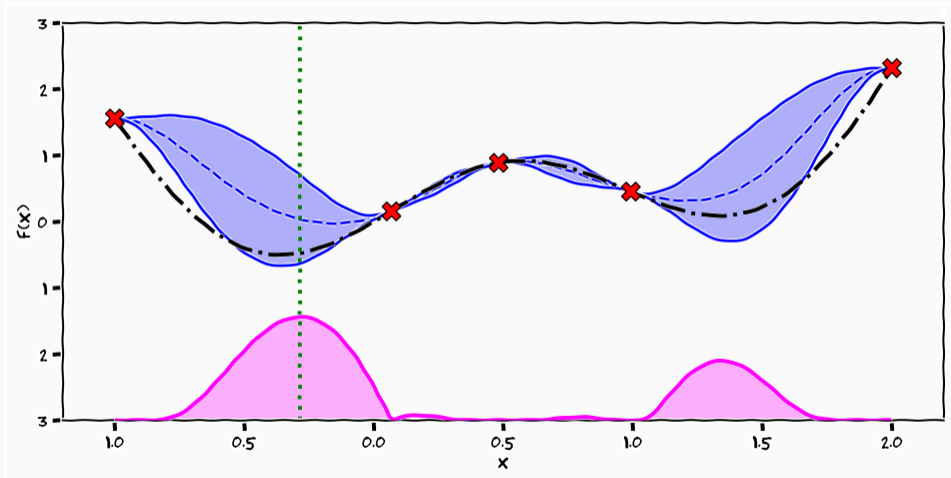


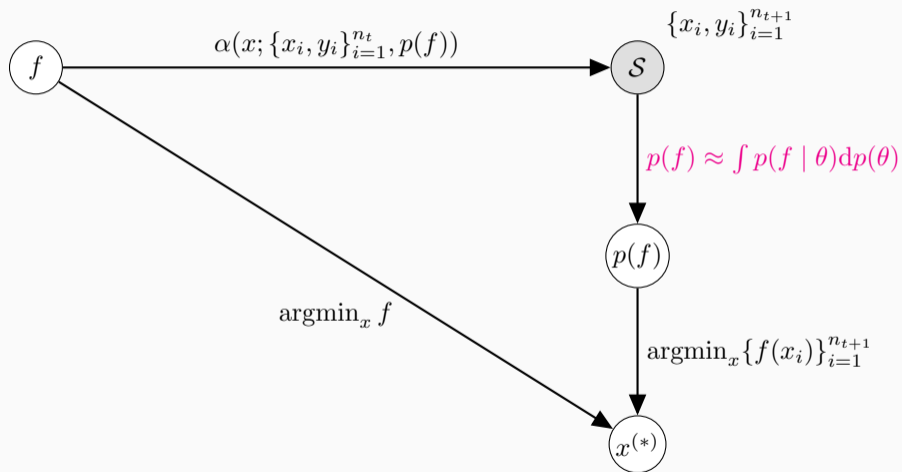
$$A \circ S \circ p(\mathcal{D}) \approx p(\mathcal{D})$$

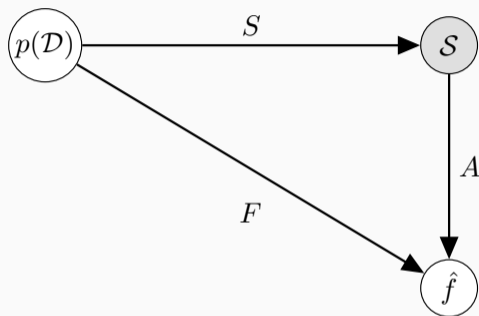


$$A \circ S \circ p(\mathcal{D}) \approx F \circ p(\mathcal{D})$$

Bayesian Optimisation







$$A \circ S \circ p(\mathcal{D}) \approx F \circ p(\mathcal{D})$$

Linear Algebra given $As = y$ estimate x s.t. $Ax = b$

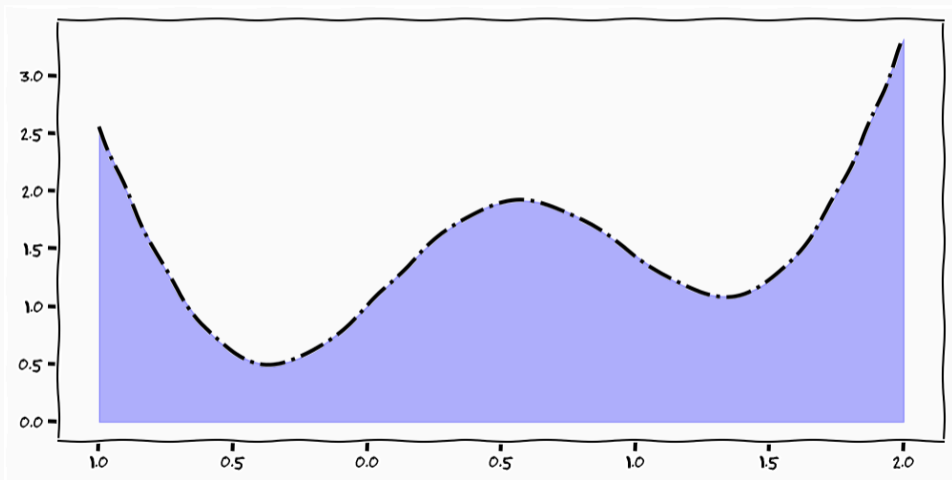
Optimisation given $\nabla f(x_i)$ estimate x s.t. $\nabla f(x) = 0$

Analysis given $f(x, t)$ estimate $x(t)$ s.t. $dx = f(x, t)$

Quadrature given $f(x_i)$ estimate $\int_a^b f(x)dx$

¹https://www.cs.toronto.edu/~duvenaud/talks/odes_runge_kutta_nips.pdf

Quantity of Interest



Integration is a significant numerical problem in many fields of science and engineering. It is a key step in inference, where it is encountered when averaging over the many states of the world consistent with observed data. Indeed, a provocative Bayesian view is that integration is the single challenge separating us from systems that fully automate statistics. More speculatively still, such systems may even exhibit artificial intelligence (ai).

– Hennig, Osbourne, Kersting

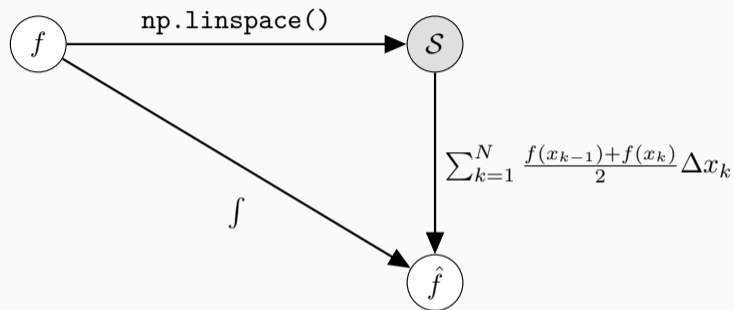
$$F := \int f(x) d\nu(x)$$

- $\nu(x)$ is the measure that we are integrating over

$$\underbrace{p(\mathcal{D})}_F = \int \underbrace{p(\mathcal{D} | \theta)}_{f(\theta)} \underbrace{p(\theta) d\theta}_{d\nu(\theta)}$$

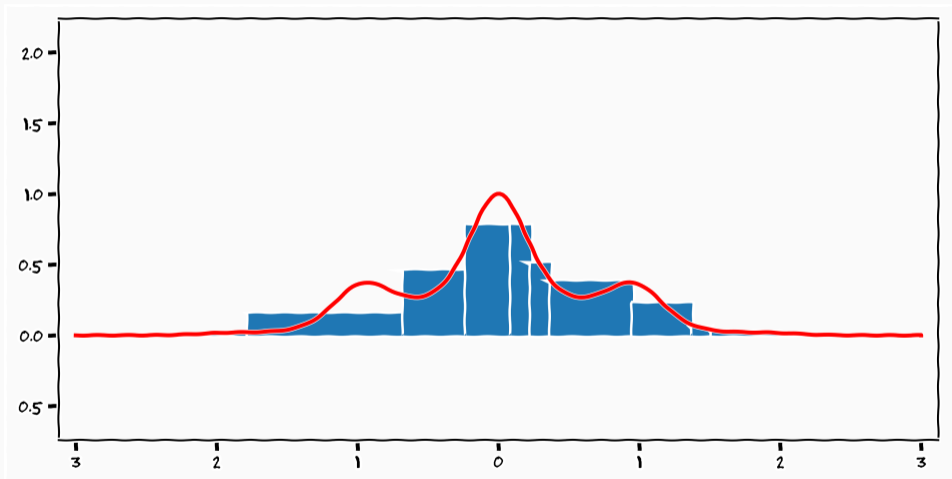
- marginalisation² is integration over the prior probability measure on the parameter

²think of computing the evidence

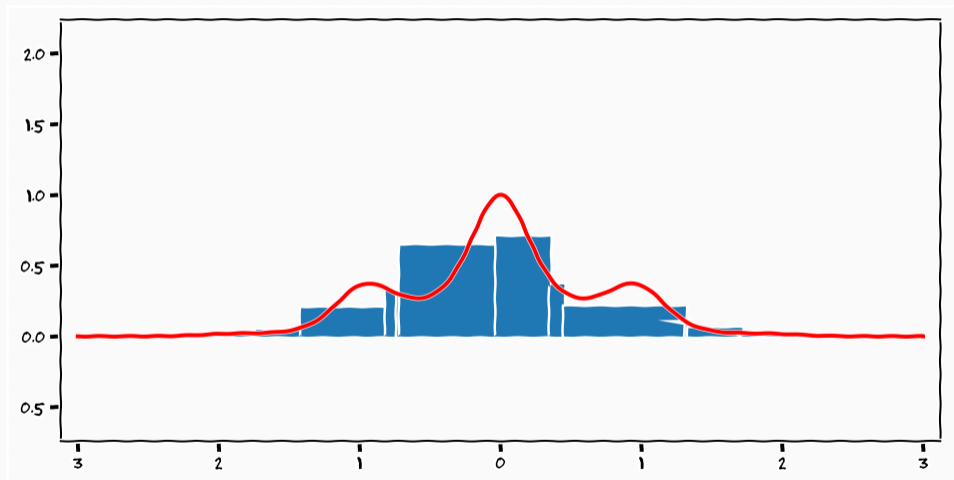


$$A \circ S \approx \int f(x) dx$$

Quadrature



Quadrature



*A numerical method **estimates** a function's **latent** property **given** the result of computations.*

A numerical method estimates a function's latent property given the result of computations.

$\frac{\text{Numerical algorithms}}{\text{Statistical inference}}$ takes data in the form of $\frac{\text{evaluations of computations}}{\text{measurements of observed variables}}$ and aims to return predictions of the quantity of interest.

A numerical method estimates a function's latent property given the result of computations.

$\frac{\text{Numerical algorithms}}{\text{Statistical inference}}$ takes data in the form of $\frac{\text{evaluations of computations}}{\text{measurements of observed variables}}$ and aims to return predictions of the quantity of interest.

Should we think about computation as inference?

$$p(\hat{f} \mid S, A)$$

Decision which algorithm to use when

$$p(\hat{f} \mid S, A)$$

Decision which algorithm to use when

Decision efficient use of expensive algorithms

$$p(\hat{f} \mid S, A)$$

Decision which algorithm to use when

Decision efficient use of expensive algorithms

Decision when to stop computation

$$p(\hat{f} \mid S, A)$$

Decision which algorithm to use when

Decision efficient use of expensive algorithms

Decision when to stop computation

Decision effect on downstream tasks

Albert Valentionvic Suldin (1924-1996) worked on error minimising estimators for numerical algorithms, how to **design** algorithms from a statistical perspective

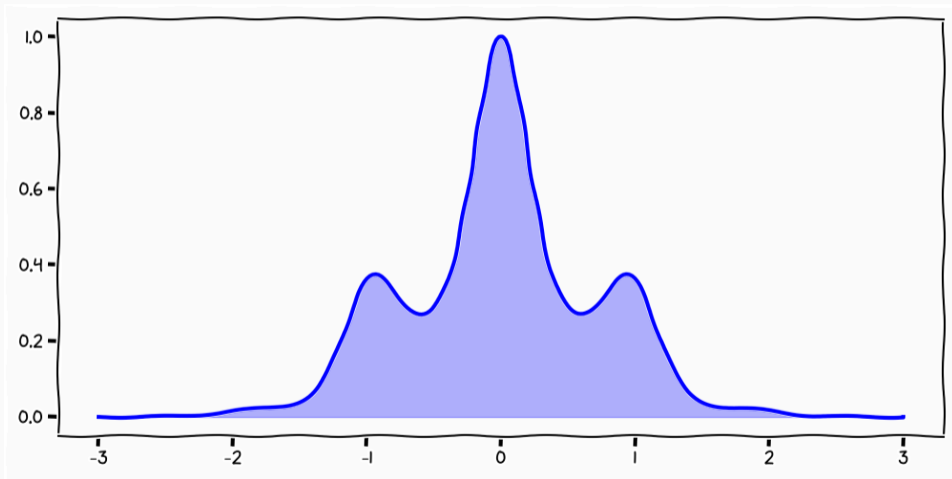
Albert Valentionvic Suldin (1924-1996) worked on error minimising estimators for numerical algorithms, how to **design** algorithms from a statistical perspective

Frederick Michael Larkin (1936-1982) incorporating the notion of **prior** knowledge into numerical algorithms to make robust calculations

Bayesian Quadrature

$$F := \int_{-3}^3 \underbrace{e^{-(\sin(3x))^2 - x^2}}_{f(x)} dx$$

- $f(x)$ fully specified and deterministic
- F is deterministic
- F cannot be computed analytically



$$p(F | Y)$$

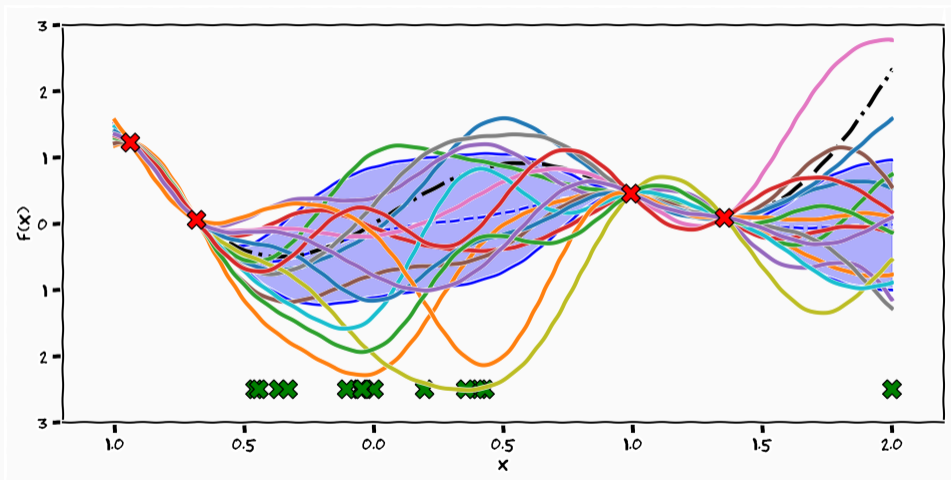
- given that I have seen data Y what is my belief about the integral

$$p(F | Y)$$

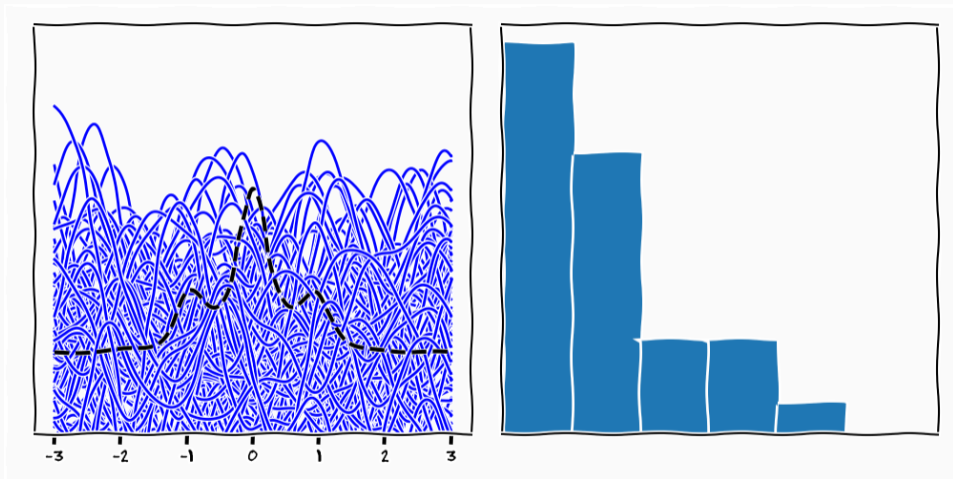
- given that I have seen data Y what is my belief about the integral
- allows for "active learning"

$$p(F | Y)$$

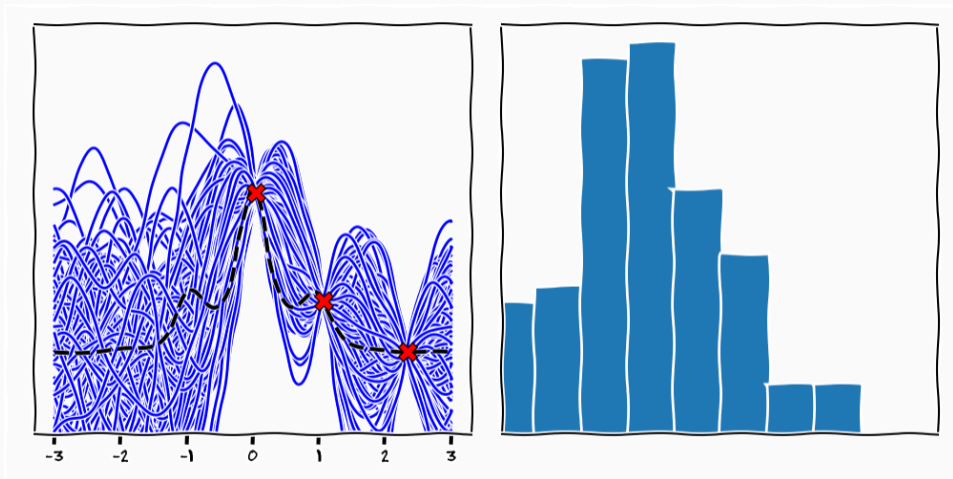
- given that I have seen data Y what is my belief about the integral
- allows for "active learning"
- exploration/exploitation etc.



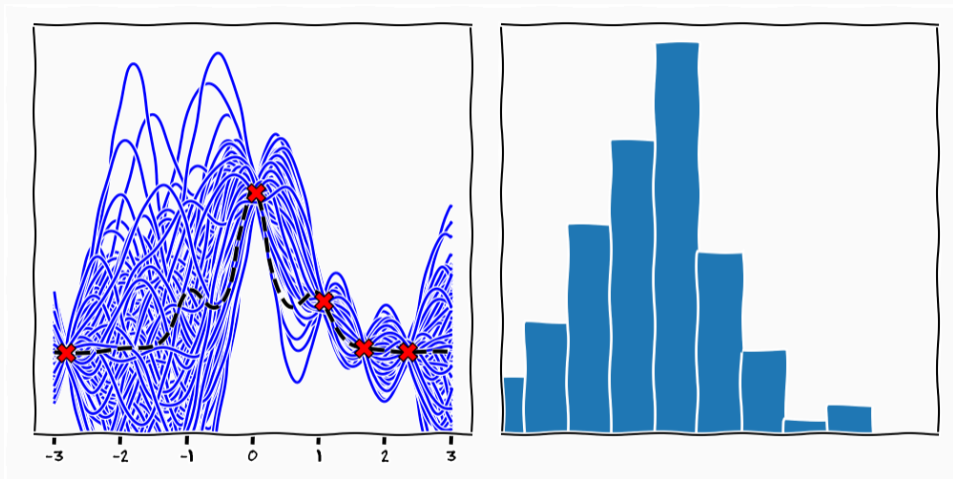
Quadrature



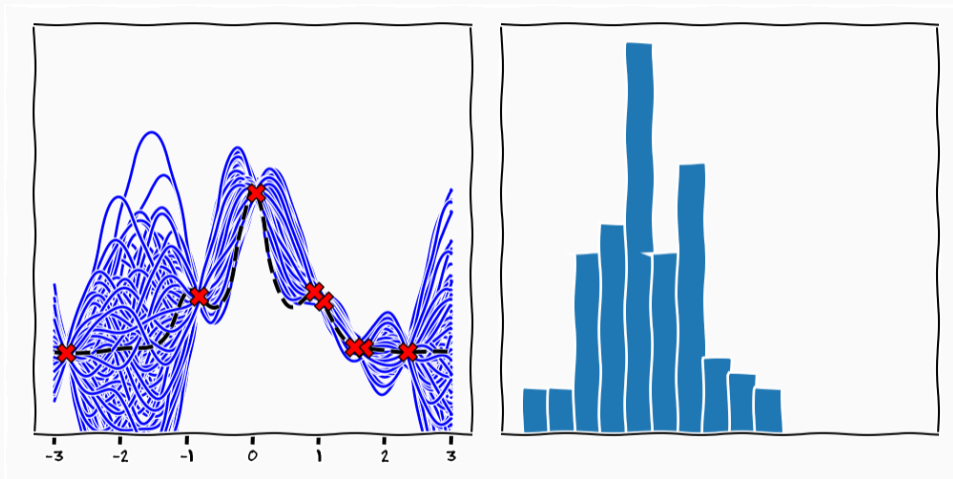
Quadrature



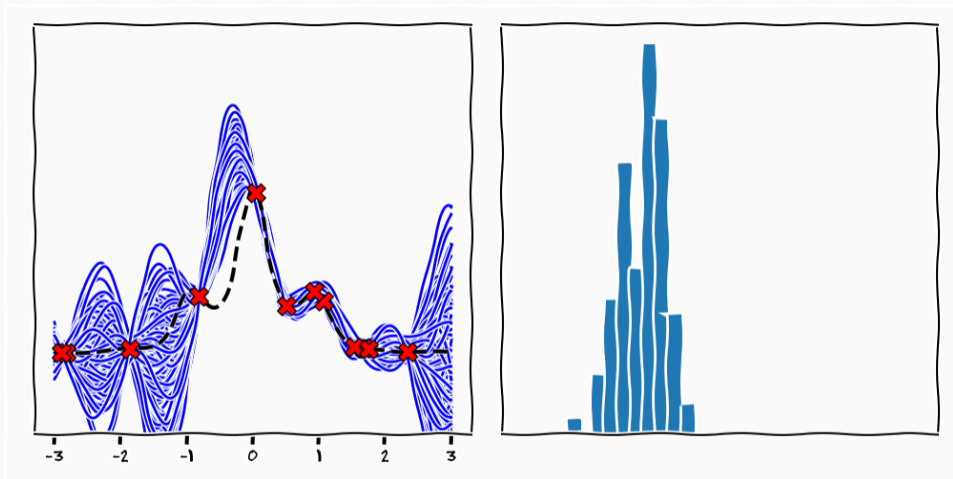
Quadrature

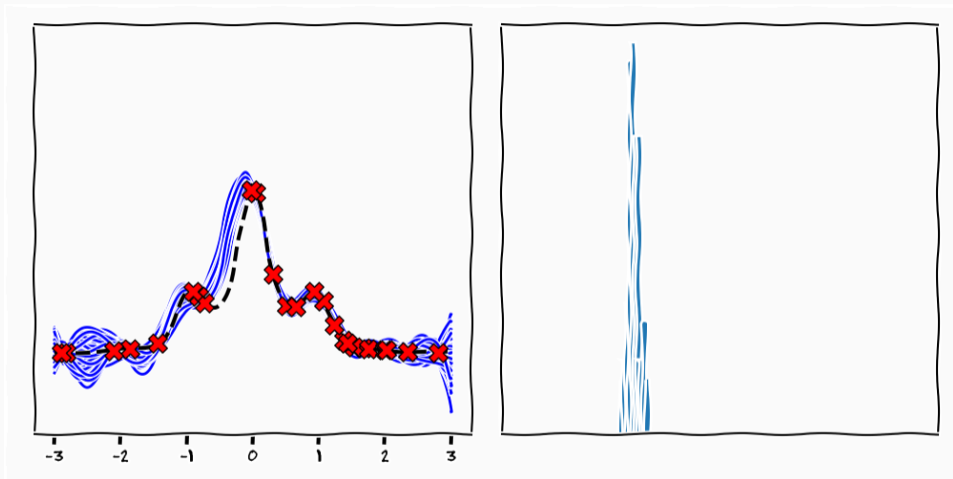


Quadrature

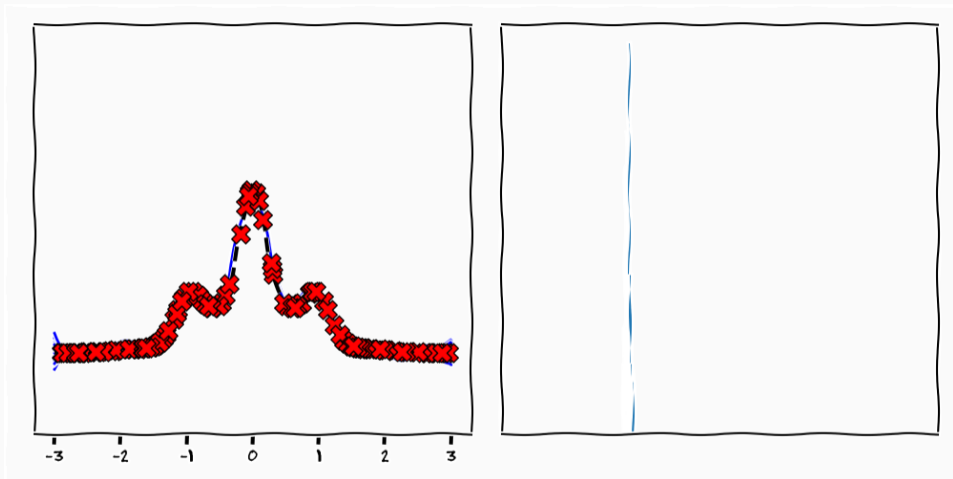


Quadrature





Quadrature



$$F := \int_{-3}^3 \underbrace{e^{-(\sin(3x))^2 - x^2}}_{f(x)} dx$$

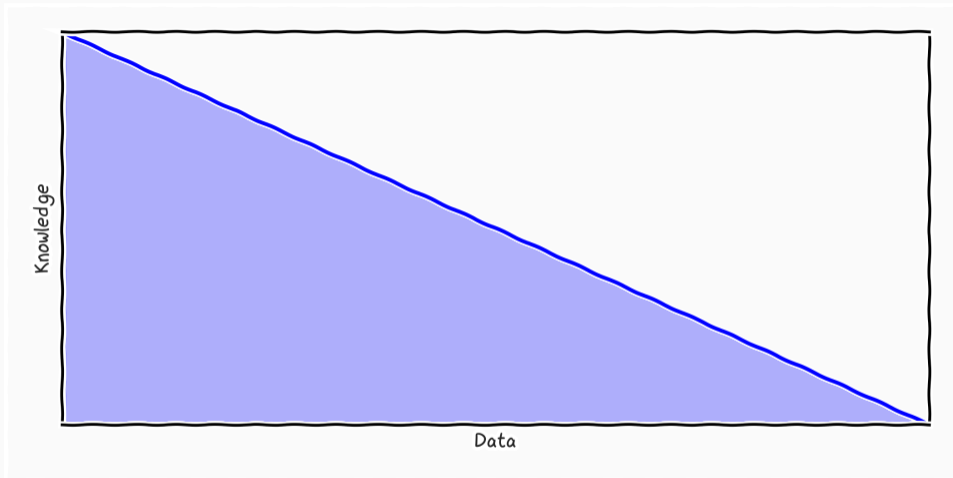
Knowledge

- $f(x)$ strictly positive $\Rightarrow F > 0$
- bounded above by,

$$f(x) \leq e^{-x^2}$$

- Therefore,

$$0 < F < \int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$



$$F := \int f(x) d\nu(x)$$

- $\nu(x)$ is the measure that we are integrating over

$$p(F, Y) = \int p(F | f)p(Y | f)p(f)df$$

$$\begin{aligned} p(F, Y) &= \int p(F | f)p(Y | f)p(f)df \\ &= \int \delta\left(F - \int_{\mathcal{X}} f dx\right) \prod_i^N \delta(y_i - f(x_i))p(f)df \end{aligned}$$

$$p \begin{pmatrix} Y \\ F \end{pmatrix} = \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_X \\ \int m_X(x) dx \end{bmatrix}, \begin{bmatrix} k(X, X) & \int k(X, x) dx \\ \int k(x, X) dx & \int \int k(x, x') dx dx' \end{bmatrix} \right)$$

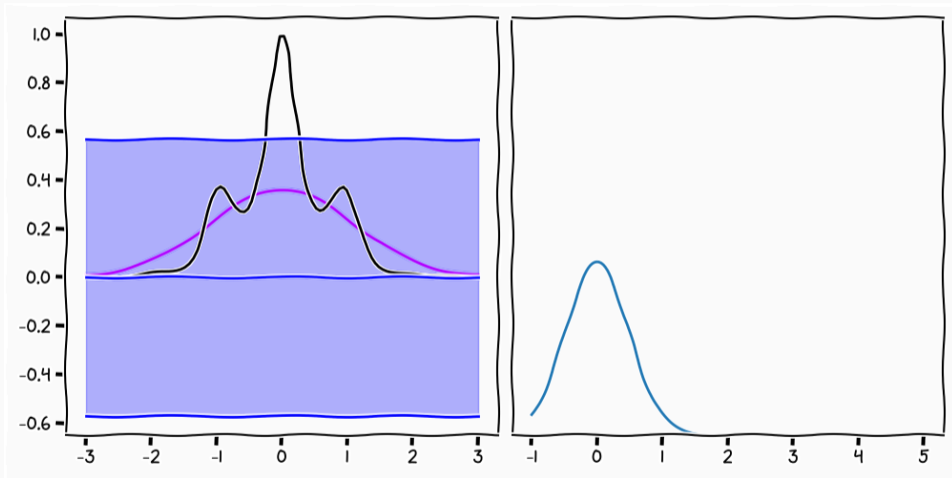
- We can derive $p(F | Y)$ through our normal conditioning procedure

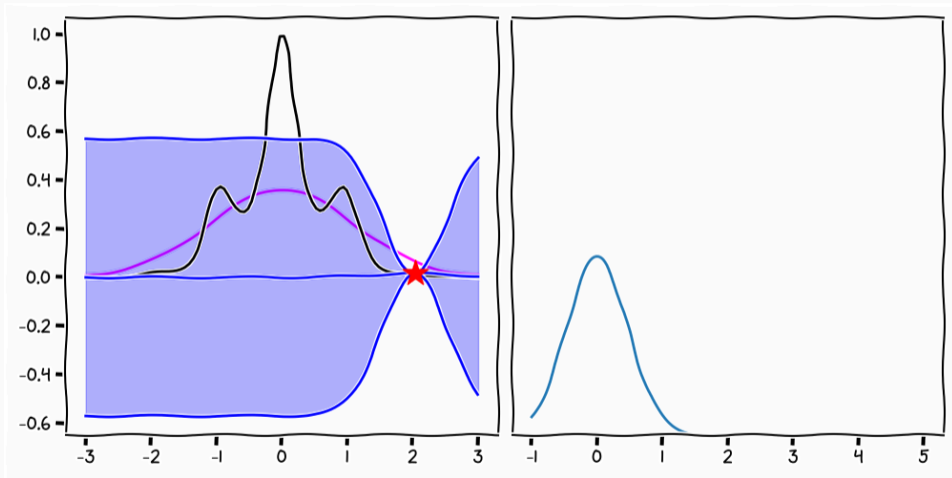
$$p \begin{pmatrix} Y \\ F \end{pmatrix} = \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_X \\ \int m_X(x) dx \end{bmatrix}, \begin{bmatrix} k(X, X) & \int k(X, x) dx \\ \int k(x, X) dx & \int \int k(x, x') dx dx' \end{bmatrix} \right)$$

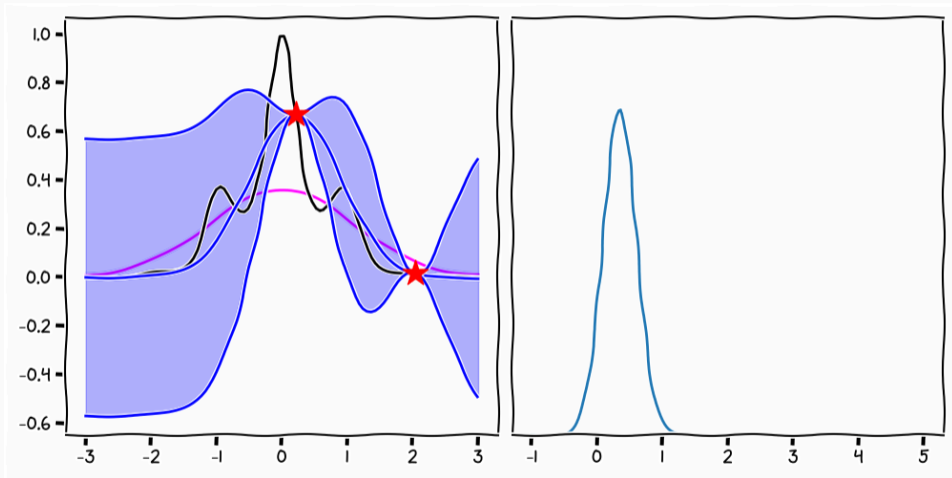
- We can derive $p(F | Y)$ through our normal conditioning procedure
- $p(F | Y) = \mathcal{N}(\mu_F, k_F)$ is a uni-variate Gaussian

$$p \begin{pmatrix} Y \\ F \end{pmatrix} = \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_X \\ \int m_X(x) dx \end{bmatrix}, \begin{bmatrix} k(X, X) & \int k(X, x) dx \\ \int k(x, X) dx & \int \int k(x, x') dx dx' \end{bmatrix} \right)$$

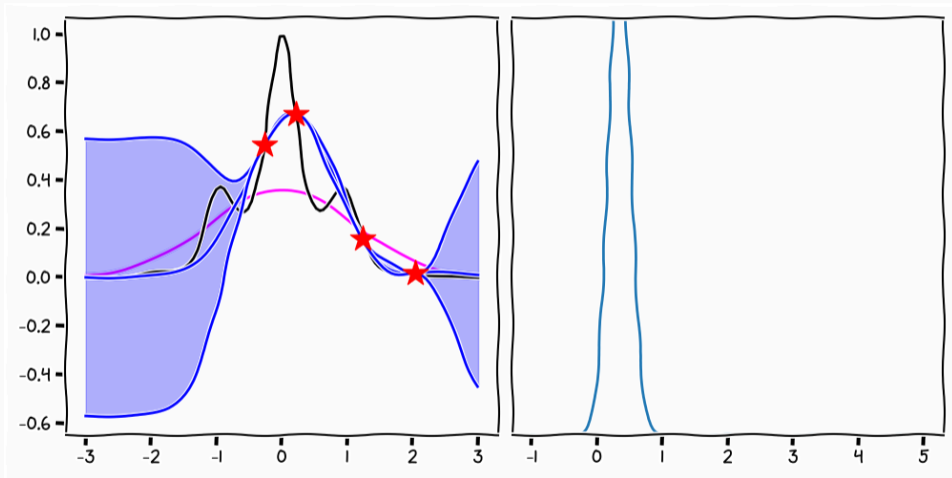
- We can derive $p(F | Y)$ through our normal conditioning procedure
- $p(F | Y) = \mathcal{N}(\mu_F, k_F)$ is a uni-variate Gaussian







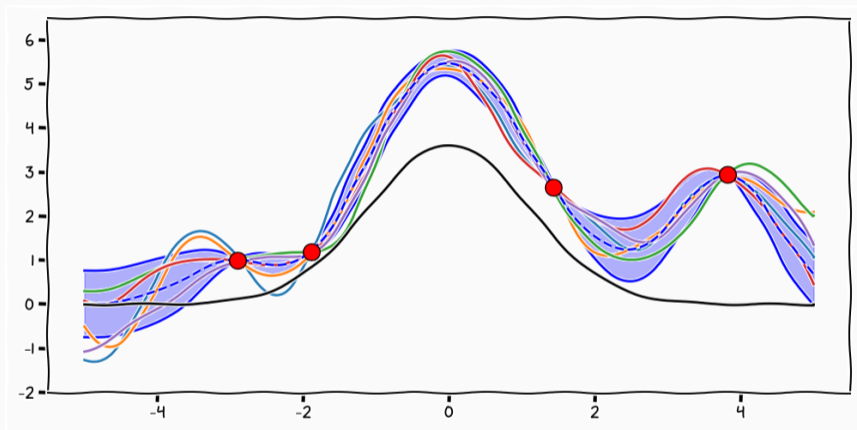
Statistical Inference



$$p \begin{pmatrix} Y \\ F \end{pmatrix} = \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_X \\ \int m_X(x) dx \end{bmatrix}, \begin{bmatrix} k(X, X) & \int k(X, x) dx \\ \int k(x, X) dx & \int \int k(x, x') dx dx' \end{bmatrix} \right)$$

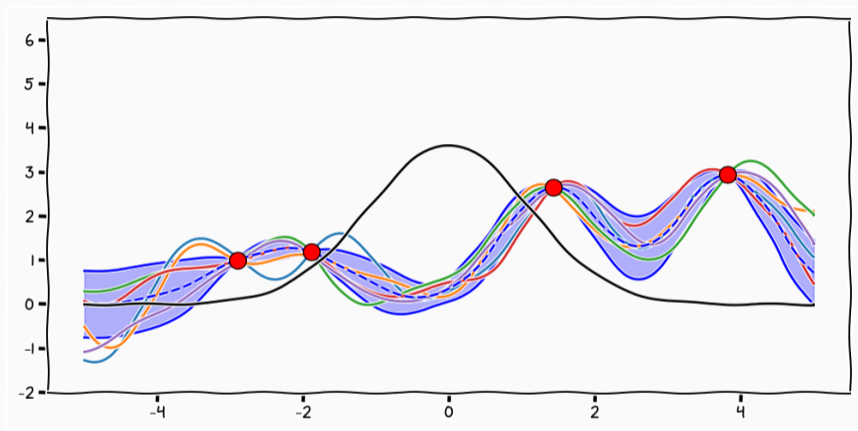
- We can derive $p(F | Y)$ through our normal conditioning procedure
- $p(F | Y) = \mathcal{N}(\mu_F, k_F)$ is a uni-variate Gaussian
- $p(Y | F) = \mathcal{N}(\mu_Y, k_Y)$ is a Gaussian process

Integral Constrained Samples

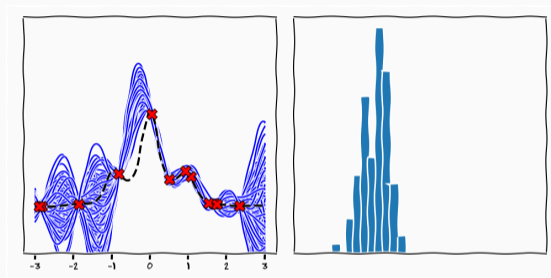


$$F = 4.0$$

Integral Constrained Samples



$$F = 1.0$$

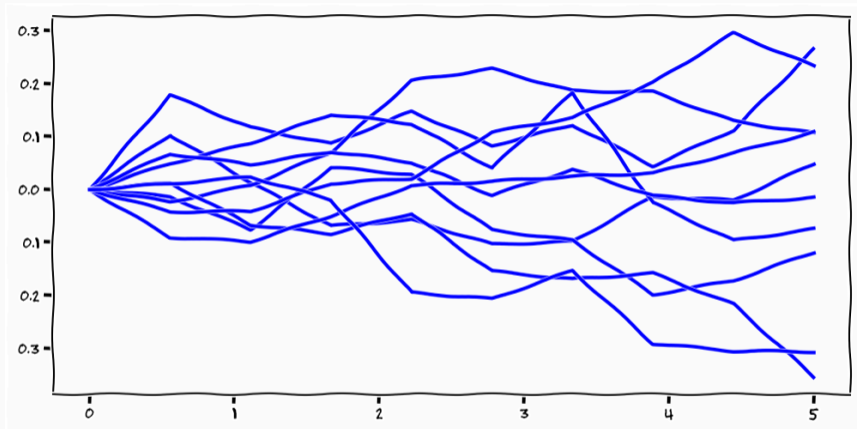


Integrand variance $\alpha(x) = k(x, x)$

Integral Variance Reduction $\alpha(x) = k_F(X, X) - k_F(X, x)$

³sometimes called a "Design Rule"

Choice of Covariance



$$p(f) = \mathcal{GP}(\mathbf{0}, \theta^2(\min(x, x') - \kappa))$$

$$\mathbb{E}[F] = \mathbb{E}_{p(f|Y)} \left[\int f(x) dx \right] = \sum_{i=1}^{N-1} \frac{x_{i+1} - x_i}{2} (f(x_{i+1}) + f(x_i))$$

$$\mathbb{E}[F] = \mathbb{E}_{p(f|Y)} \left[\int f(x) dx \right] = \sum_{i=1}^{N-1} \frac{x_{i+1} - x_i}{2} (f(x_{i+1}) + f(x_i))$$

- This is the normal trapezoid rule!!!

$$\mathbb{E}[F] = \mathbb{E}_{p(f|Y)} \left[\int f(x) dx \right] = \sum_{i=1}^{N-1} \frac{x_{i+1} - x_i}{2} (f(x_{i+1}) + f(x_i))$$

- This is the normal trapezoid rule!!!
- The algorithm is now tied to **our belief** in the function!!!!

$$\mathbb{E}[F] = \mathbb{E}_{p(f|Y)} \left[\int f(x) dx \right] = \sum_{i=1}^{N-1} \frac{x_{i+1} - x_i}{2} (f(x_{i+1}) + f(x_i))$$

- This is the normal trapezoid rule!!!
- The algorithm is now tied to **our belief** in the function!!!!
- We can do inference over where to sample!!!!!!!!!!

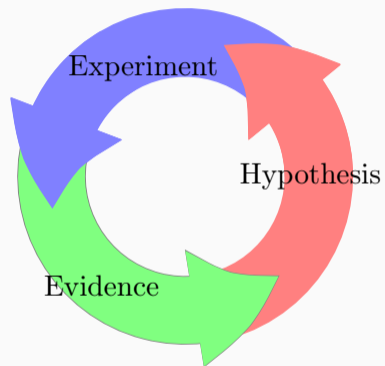
Definition (Trapezoid Rule)

The trapezoidal rule is the posterior mean estimate for the integral

$F = \int_a^b f(x)dx$ under any centred Wiener process prior $p(f) = \mathcal{GP}(0, k)$ with $k(x, x') = \theta^2(\min(x, x') - \kappa)$

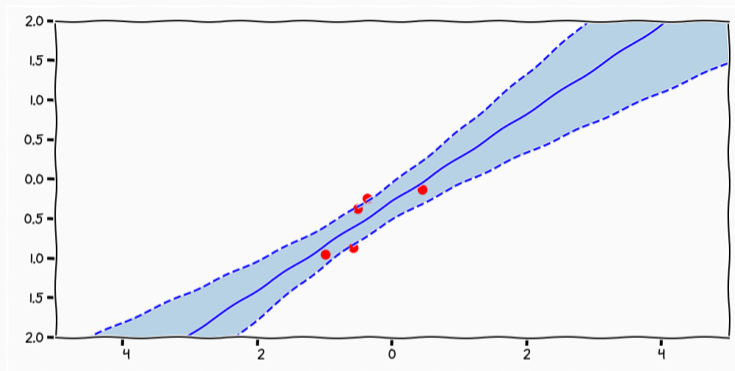
I'M NOT IMPRESSED





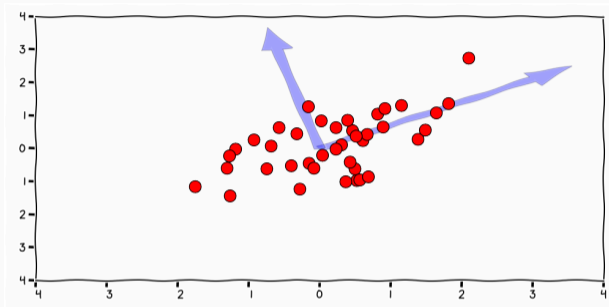
Compute
Data + Model $\xrightarrow{\quad}$ Prediction

Least Squares Regression



Legendre (1805) algorithm that reduces "error"

Gauss (1809) statistical model assuming i.i.d. Gaussian noise



Spearman (1904) proposed an algorithm to extract "factors" from data
Spearman, [1904](#)

Hotelling (1936) concept of factor **is** clearly defined through a statistical
model Hotelling, [1933](#)

Code

```
def minimize(fun, x0, args=(), method=None,
            jac=None, hess=None,
            hessp=None, bounds=None,
            constraints=(), tol=None,
            callback=None, options=None):
```

method Nelder-Mead, Powell, CG, BFGS, Newton-CG, L-BFGS-B, TNC ,
COBYLA , SLSQP , trust-constr , dogleg , trust-ncg ,
trust-exact , trust-krylov

- There are tons of numerical algorithms for every problem under the sun

- There are tons of numerical algorithms for every problem under the sun
- They work really well

- There are tons of numerical algorithms for every problem under the sun
- They work really well
- They give different results on the same problem

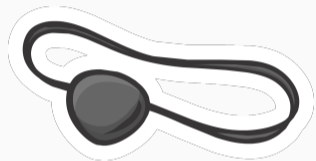
- There are tons of numerical algorithms for every problem under the sun
- They work really well
- They give different results on the same problem
- *what is the prior they implement?*

In a talk, Olivier Bousquet has described the deep learning community as a giant genetic algorithm: Researchers in this community are exploring the space of all variants of algorithms and architectures in a semi-random way. Things that consistently work in large experiments are kept, the ones not working are discarded. the community is evolving only one set of parameters (architectures, initialization strategies, hyperparameters search algorithms, etc.) keeping most of the time the optimizer fixed to Adam.

"There is a notion of success . . . which I think is novel in the history of science. It interprets success as approximating unanalyzed data."
– Prof. Noam Chomsky

⁴Chomsky et al., 1980





Statistical Learning



$$A_{\mathcal{H}}(S)$$



Summary

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵

⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties

⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties
- Computation is often not "truth" therefore we should quantify our ignorance about its results

⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties
- Computation is often not "truth" therefore we should quantify our ignorance about its results
- Why?

⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties
- Computation is often not "truth" therefore we should quantify our ignorance about its results
- Why?
 - efficiency

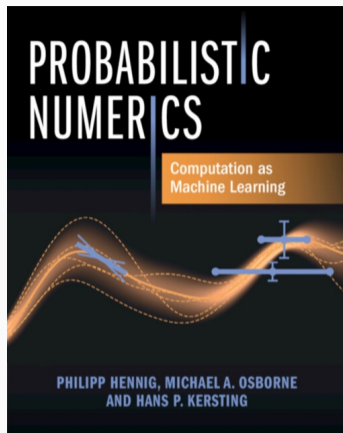
⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties
- Computation is often not "truth" therefore we should quantify our ignorance about its results
- Why?
 - efficiency
 - down-stream tasks, uncertainty in computation should be part of decision

⁵these thoughts have been around for a long time

- Probabilistic Numerics extends the notion of statistical inference to **computation**⁵
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties
- Computation is often not "truth" therefore we should quantify our ignorance about its results
- Why?
 - efficiency
 - down-stream tasks, uncertainty in computation should be part of decision
 - learning/understanding algorithms in relation to problems/data

⁵these thoughts have been around for a long time



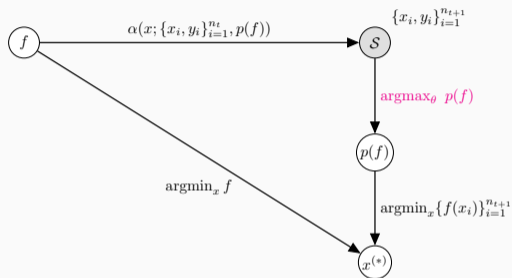
<http://probnumschool.org>

Adaptive probabilistic ODE solvers without adaptive memory requirements

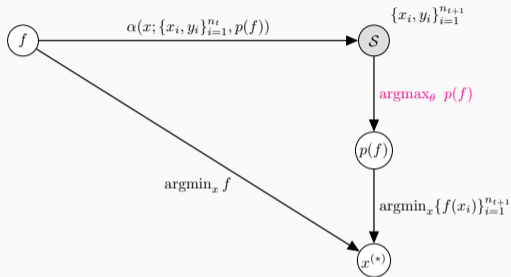
Adaptive probabilistic solvers for ordinary differential equations (ODEs) have made substantial progress in recent years but can still not solve memory-demanding differential equations. In this talk, I review recent developments in numerically robust fixed-point smoothers and how to use them for constructing adaptive probabilistic ODE solvers. These new algorithms use drastically less memory than their predecessors and are the first adaptive probabilistic numerical methods compatible with scientific computing in JAX .

SS03, 16-17

eof

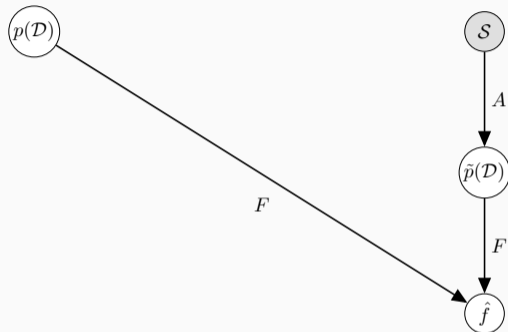


Yes it uses a probabilistic model as a proxy for decision loop






Yes it uses a probabilistic model as a proxy for decision loop





No the probabilistic model is not over the quantity of interest



$$A \circ \mathcal{S} = \tilde{p}(\mathcal{D}) \approx p(\mathcal{D})$$

References

-  Chomsky, Noam A and Jerry A Fodor (1980). **“The inductivist fallacy.”** In: *Language and Learning: The Debate between Jean Piaget and Noam Chomsky*.
-  Cockayne, Jon, Chris Oates, Tim Sullivan, and Mark Girolami (2017). **“Bayesian Probabilistic Numerical Methods.”** In: *CoRR*.
-  Hennig, Philipp, Michael A Osborne, and Mark Girolami (July 2015). **“Probabilistic numerics and uncertainty in computations.”** In: *Proc. R. Soc. A* 471.2179, p. 20150142.

-  Hotelling, H (Sept. 1933). “**Analysis of a complex of statistical variables into principal components..**” In: *Journal of Educational Psychology* 24.6, pp. 417–441.
-  Neumann, John von and H. H. Goldstine (1947). “**Numerical Inverting of Matrices of High Order.**” In: *Bulletin of the American Mathematical Society* 53.11, pp. 1021–1100.
-  O’Hagan, A. (Nov. 1991). “**Bayes-Hermite quadrature.**” In: *Journal of Statistical Planning and Inference* 29.3, pp. 245–260.
-  Spearman, Charles (1904). “**" General Intelligence," Objectively Determined and Measured.**” In: *The American Journal of Psychology* 15.2, pp. 201–292.