



UNIVERSITY OF
CAMBRIDGE

Machine Learning and the Physical World

Lecture 7 : Probabilistic Numerics

Carl Henrik Ek - che29@cam.ac.uk

31st of October, 2023

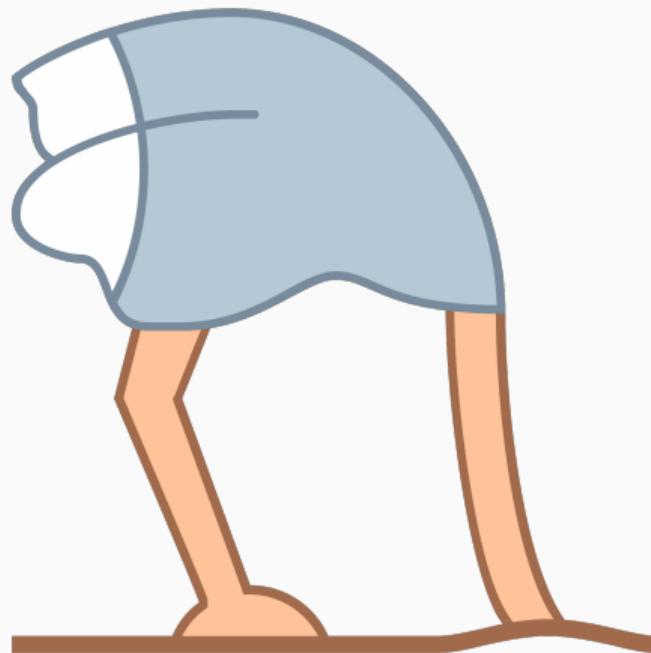
<http://carlhenrik.com>

Compute
Data + Model $\overbrace{\rightarrow}^{\text{Compute}}$ Prediction

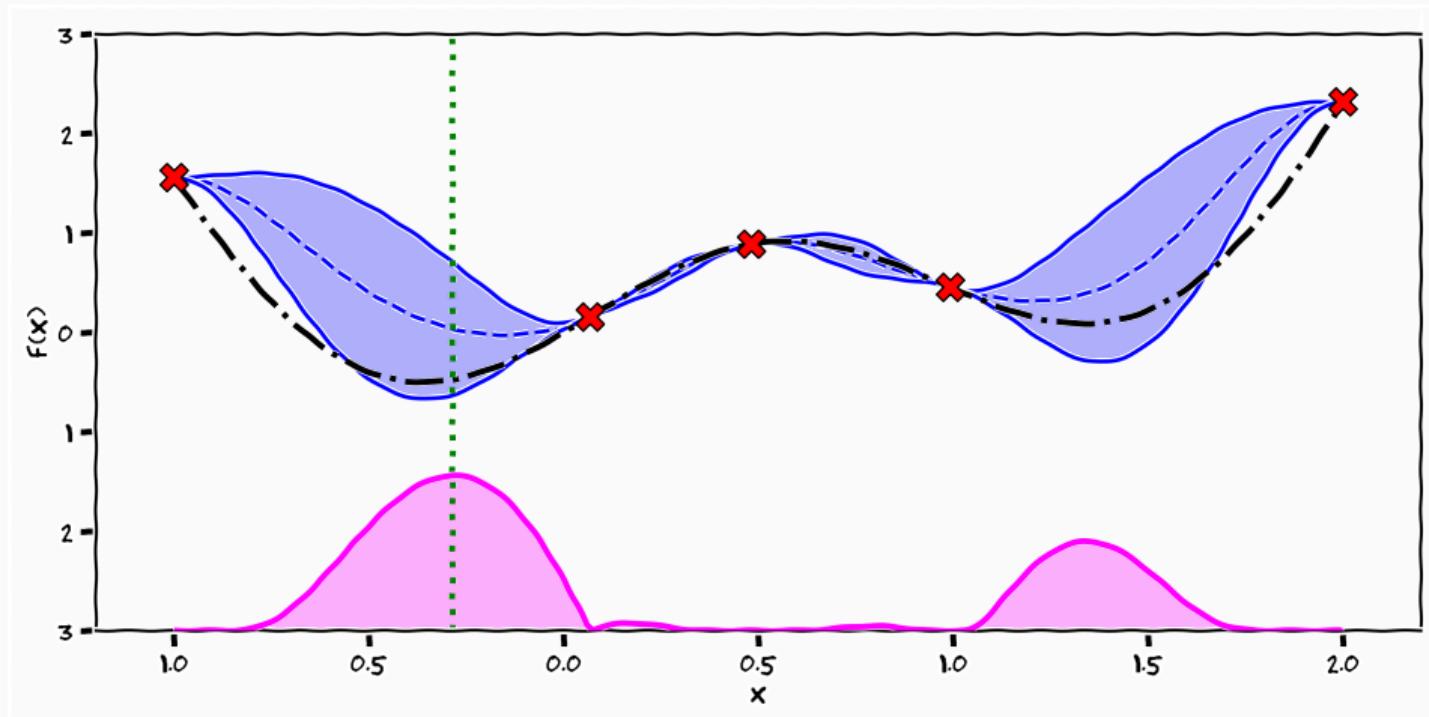
The role of Uncertainty/Ignorance

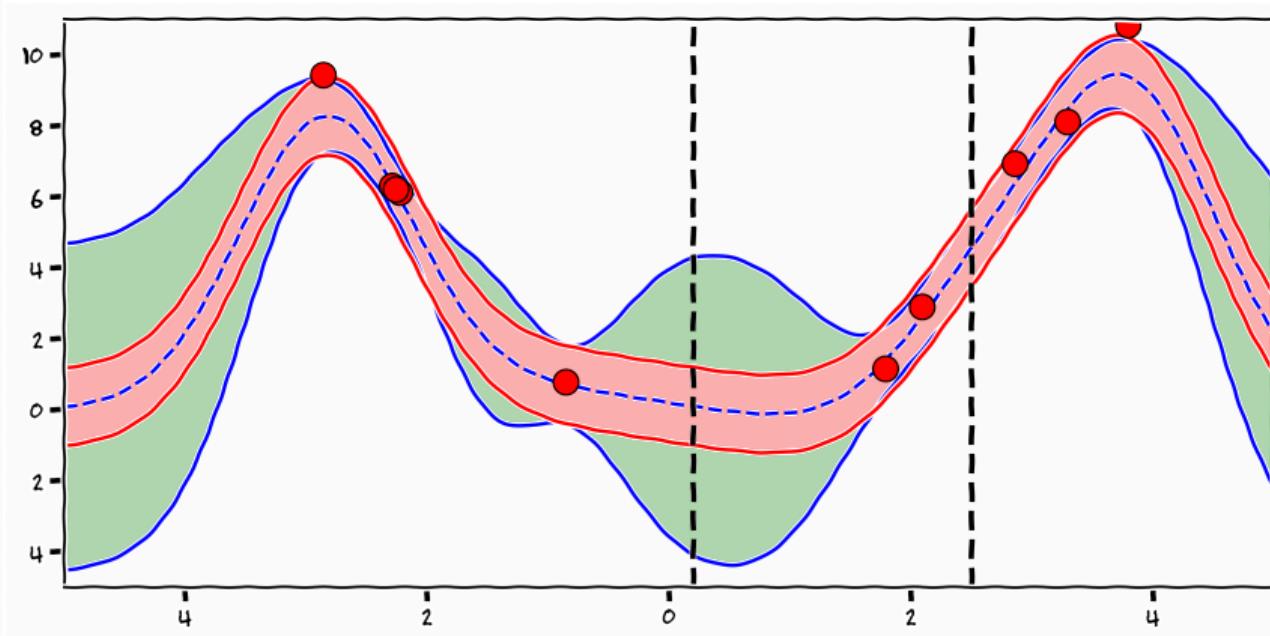
$$p(y) = \int p(y \mid f)p(f)df$$

What do we do with uncertainty?



Bayesian Optimisation





$$y_i = f_i + \epsilon$$

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 the underlying system

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

Epistemic Uncertainty related to our ignorance of the underlying system

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

"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 Poincaré, 1896

Computational Decisions

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

Computational Decisions

- 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?

Computational Decisions

- 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?

Computational Decisions

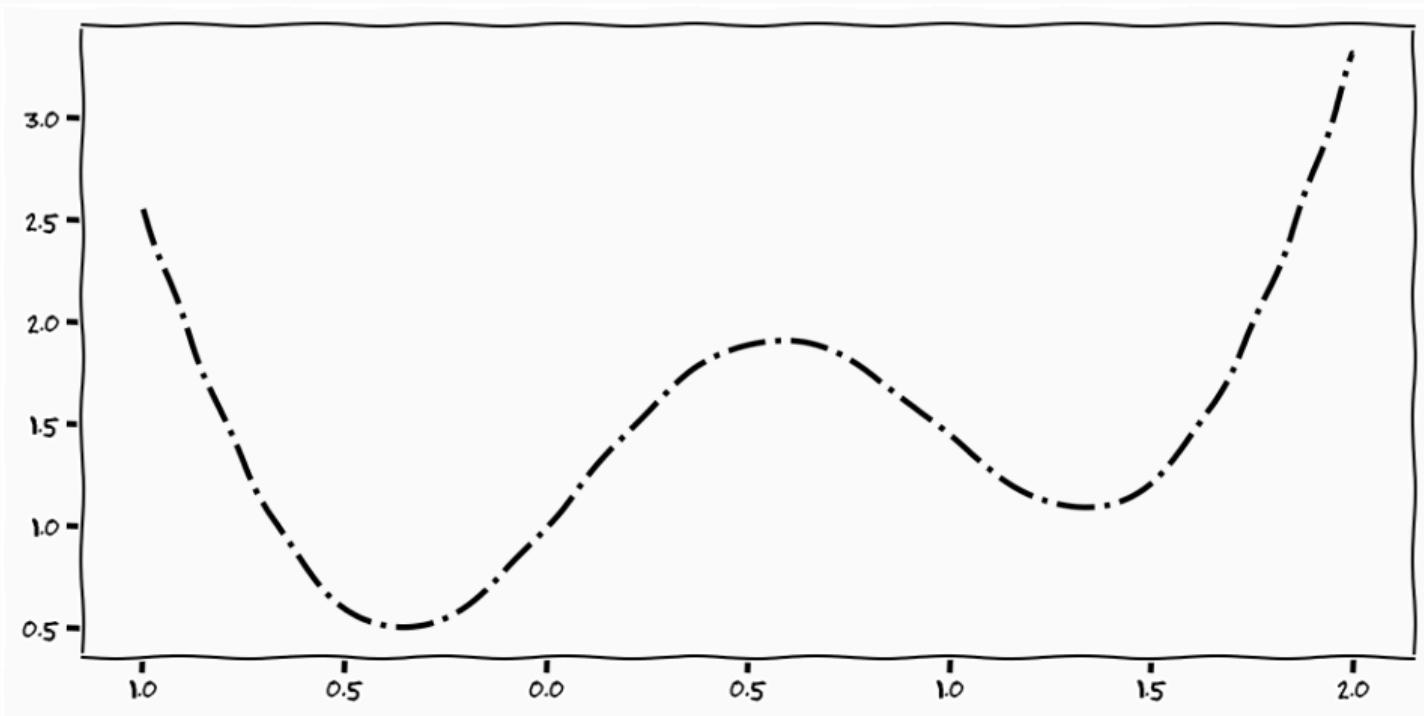
- Computation is expensive, how much knowledge will I gain from computing more?
- What should I compute in order to reduce my uncertainty as much as possible?
- How much should I trust the computation I have done?
- How precise should I do down-stream tasks based on the information from a specific computation?

Why Probabilistic Numerics?

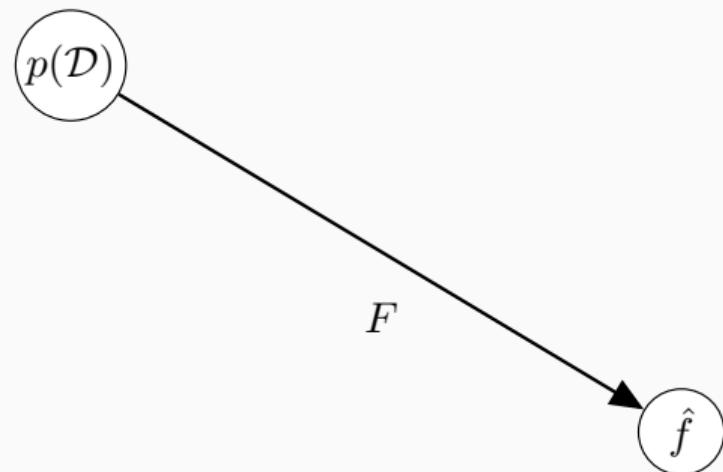
"[round-off errors] are strictly very complicated but uniquely defined number theoretical functions [of the inputs], yet our ignorance of their true nature is such that we best treat them as random variables."

– Neumann et al., 1947

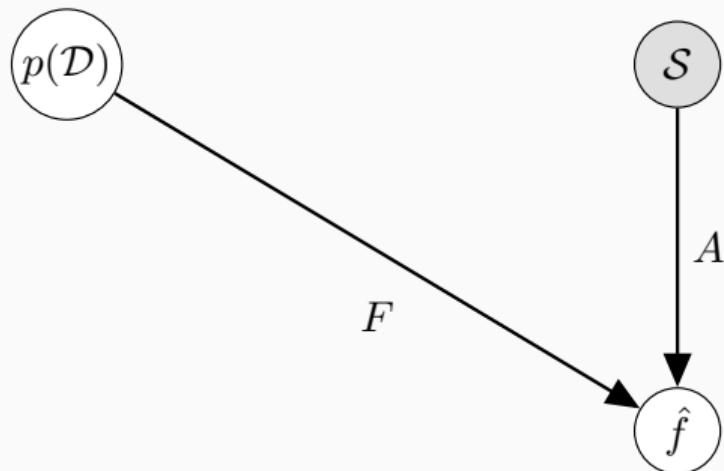
I believe in. . .



Formalisation [Cockayne et al., 2017]

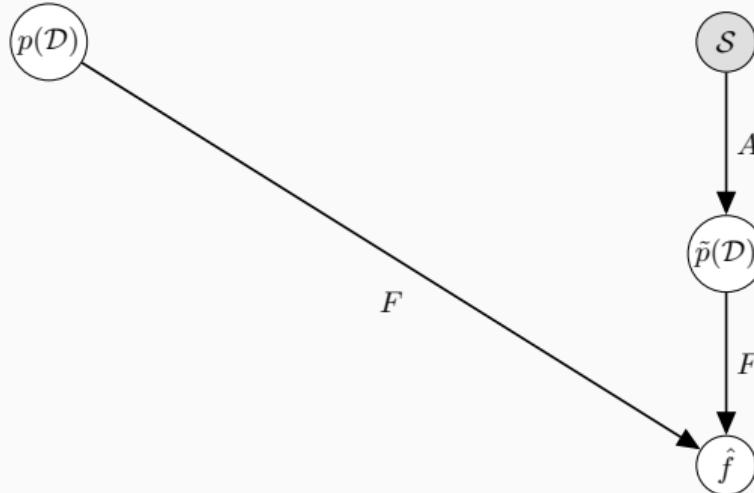


$$F : p(\mathcal{D}) \rightarrow p(\mathcal{Y}|\mathcal{X})$$

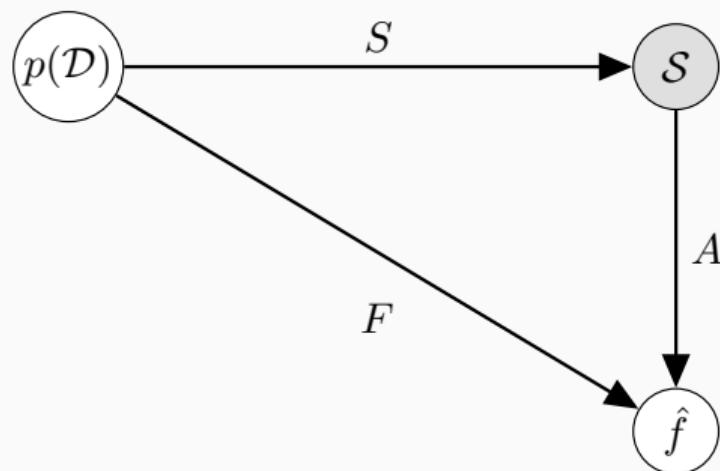


$$A \circ \mathcal{S} \approx F \circ p(\mathcal{D})$$

Formalisation

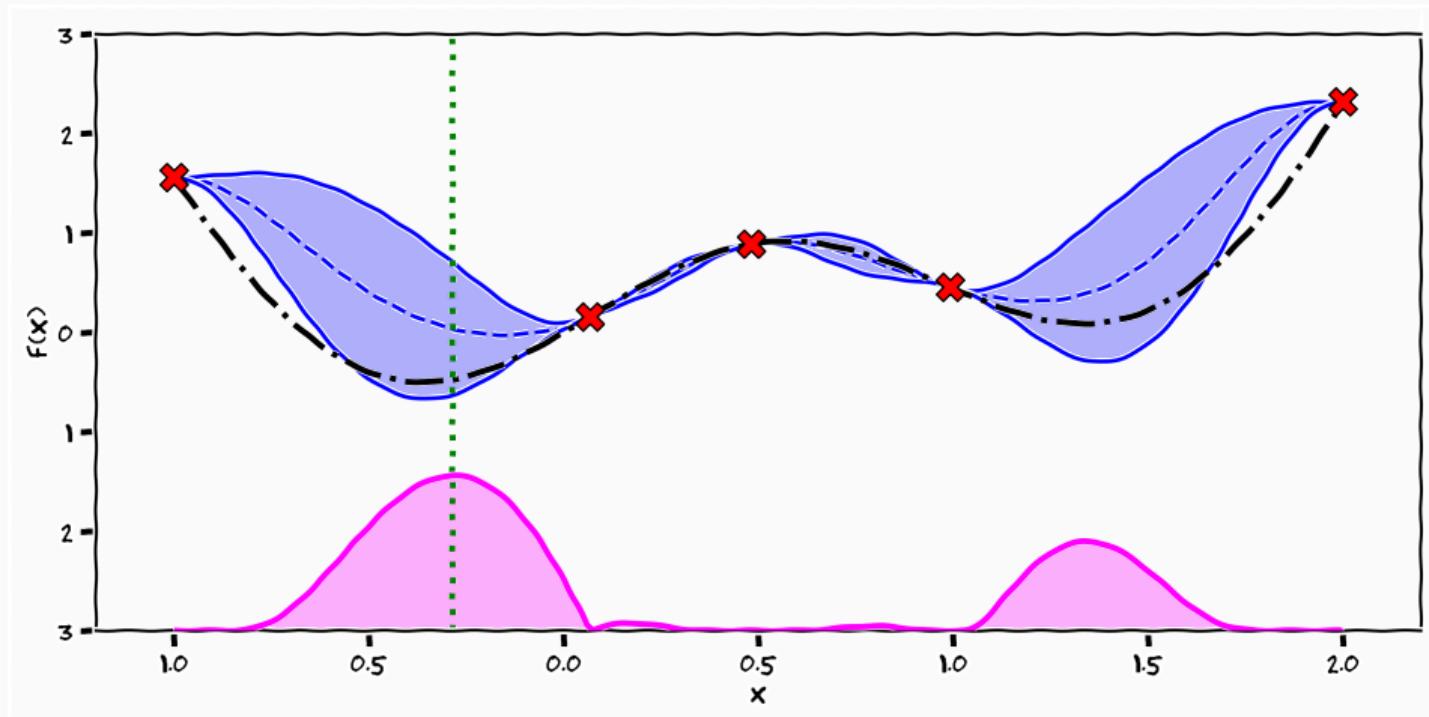


$$A \circ S \circ p(\mathcal{D}) \approx p(\mathcal{D})$$

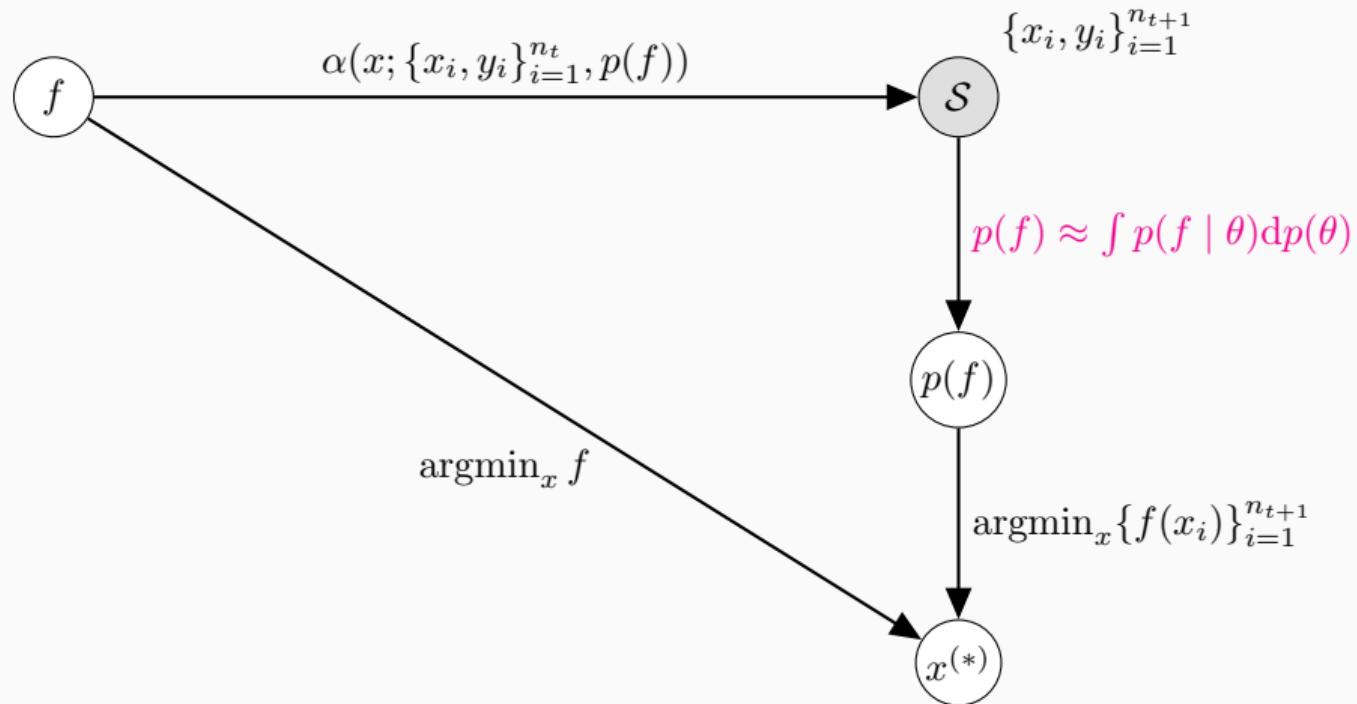


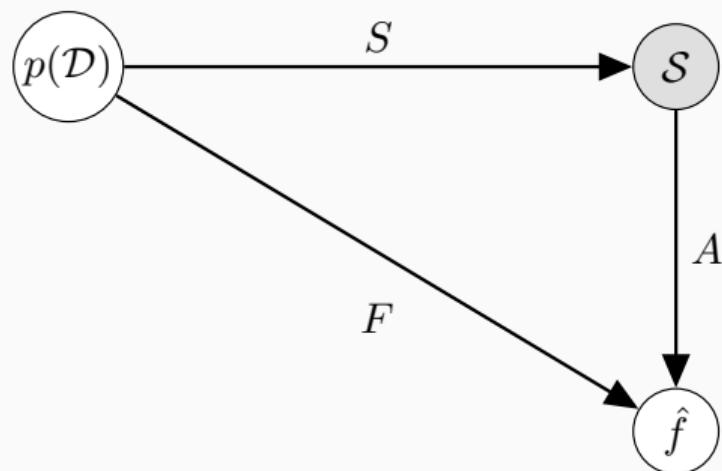
$$A \circ S \circ p(\mathcal{D}) \approx F \circ p(\mathcal{D})$$

Bayesian Optimisation



Bayesian Optimisation





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

Numerical Computations¹

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

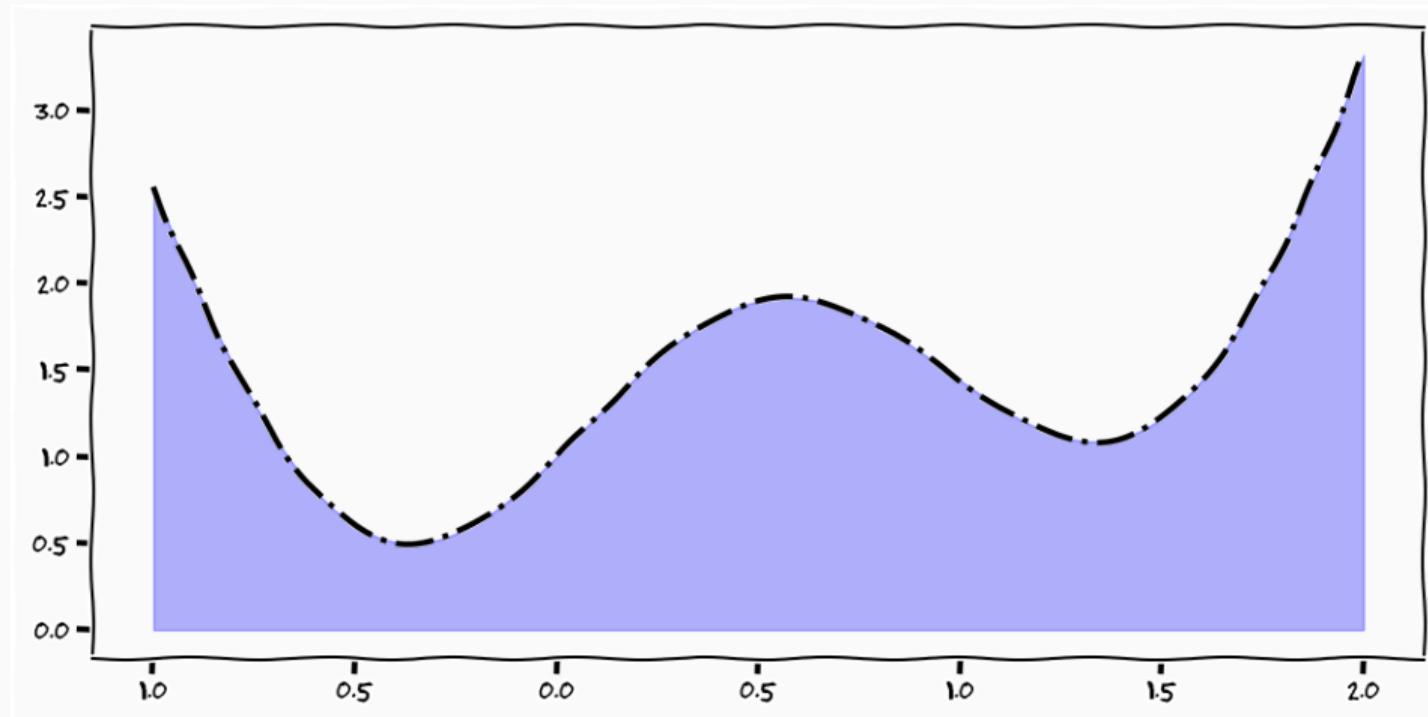
Optimisation given $\nabla f(x_i)$ estimate x s.t. $\nabla f(x) = 0$

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

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

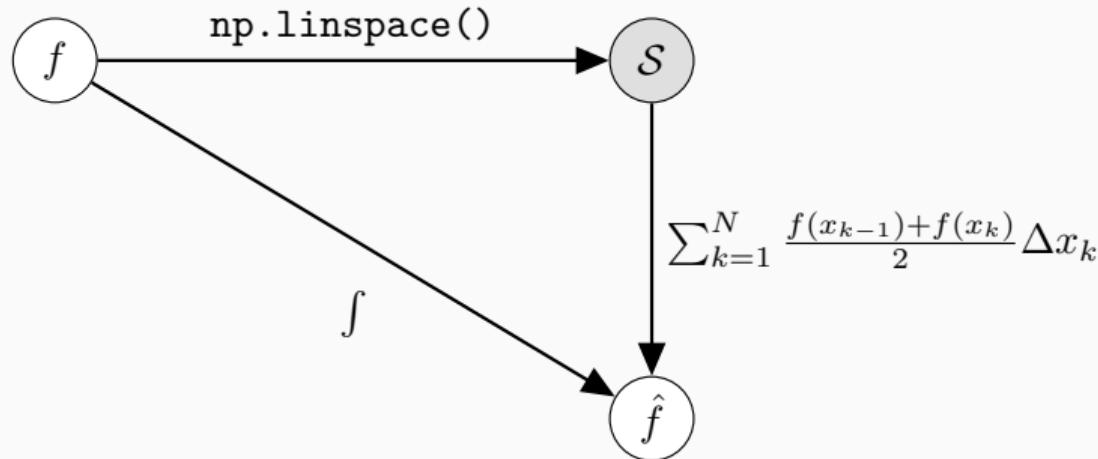
- $\nu(x)$ is the measure that we are integrating over

Integration

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

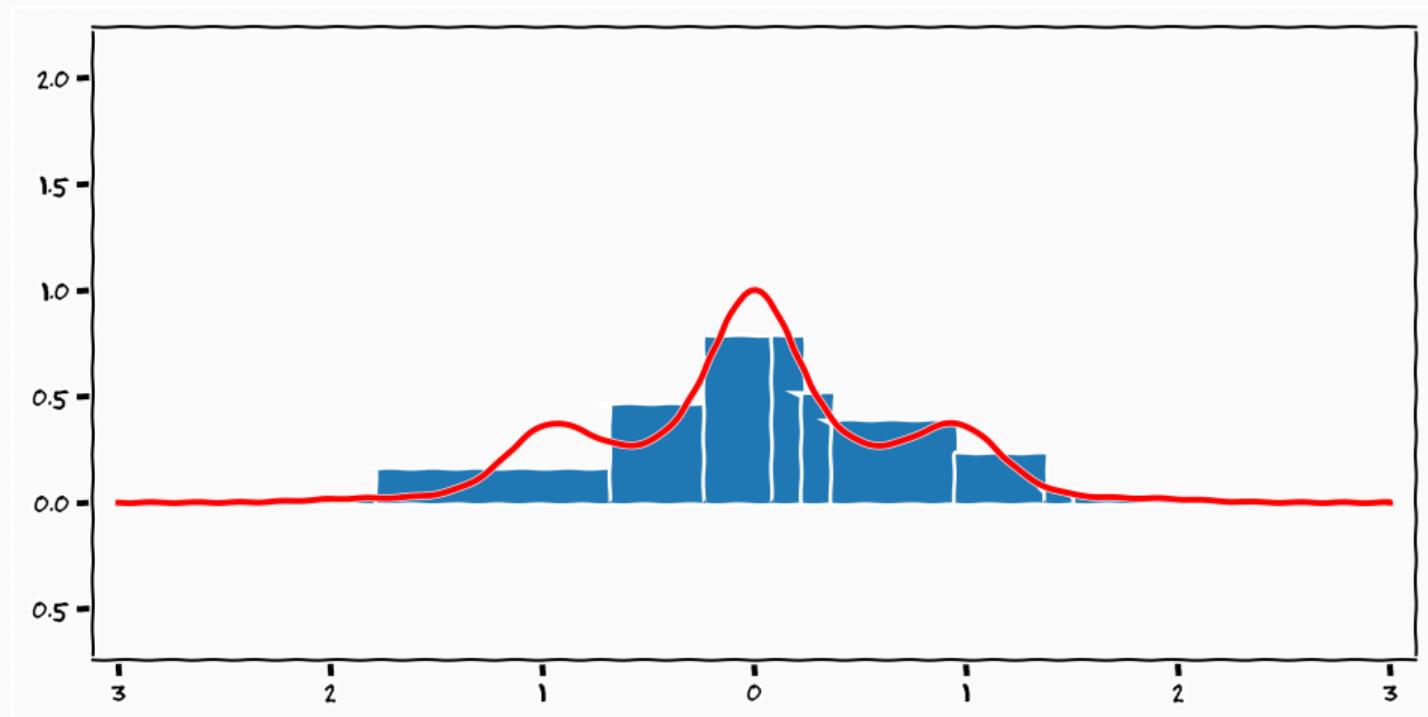
- marginalisation² is integration over the prior probability measure on the parameter

²think of computing the evidence

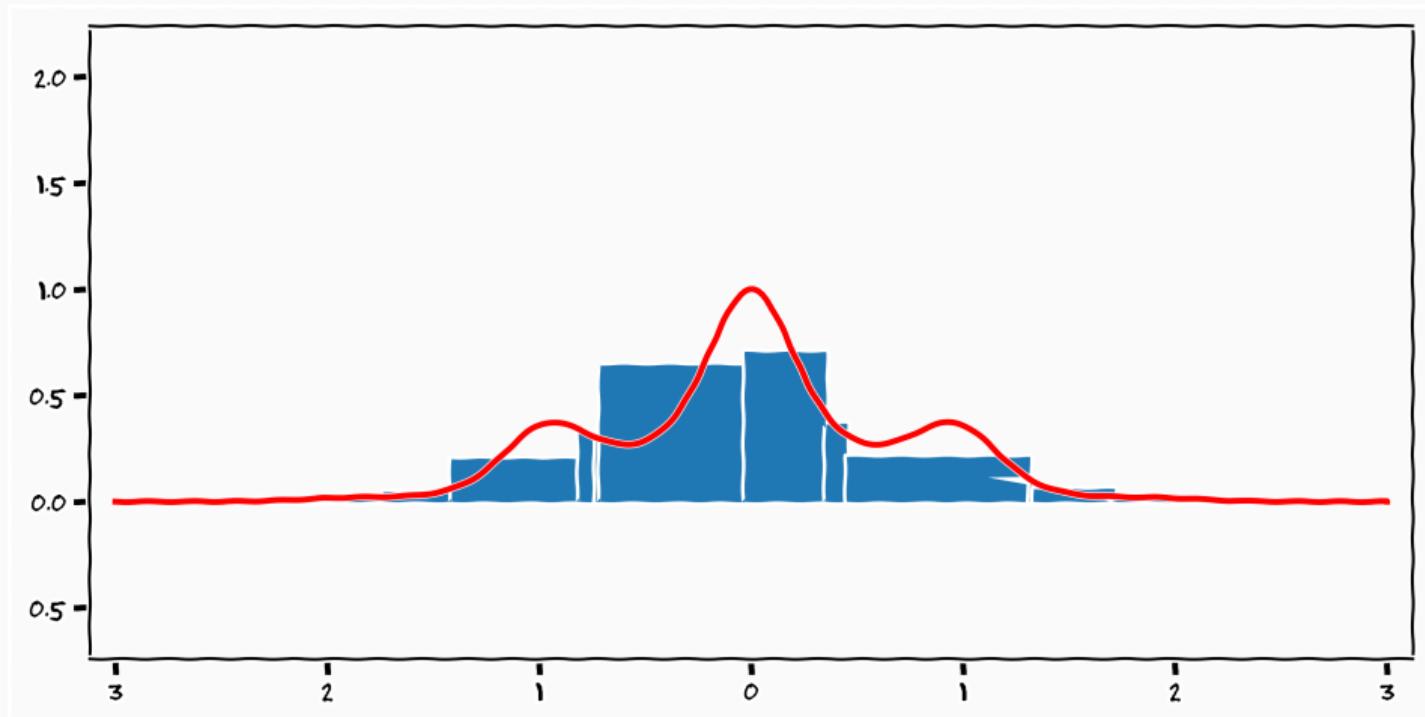


$$A \circ \mathcal{S} \approx \int f(x) dx$$

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.

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

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

Numerical algorithms
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?

Use of Computational Uncertainty

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

Decision which algorithm to use when

Use of Computational Uncertainty

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

Decision which algorithm to use when

Decision efficient use of expensive algorithms

Use of Computational Uncertainty

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

Decision which algorithm to use when

Decision efficient use of expensive algorithms

Decision when to stop computation

Use of Computational Uncertainty

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

Decision which algorithm to use when

Decision efficient use of expensive algorithms

Decision when to stop computation

Decision effect on downstream tasks

When computation was expensive

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

When computation was expensive

Albert Valentionic 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

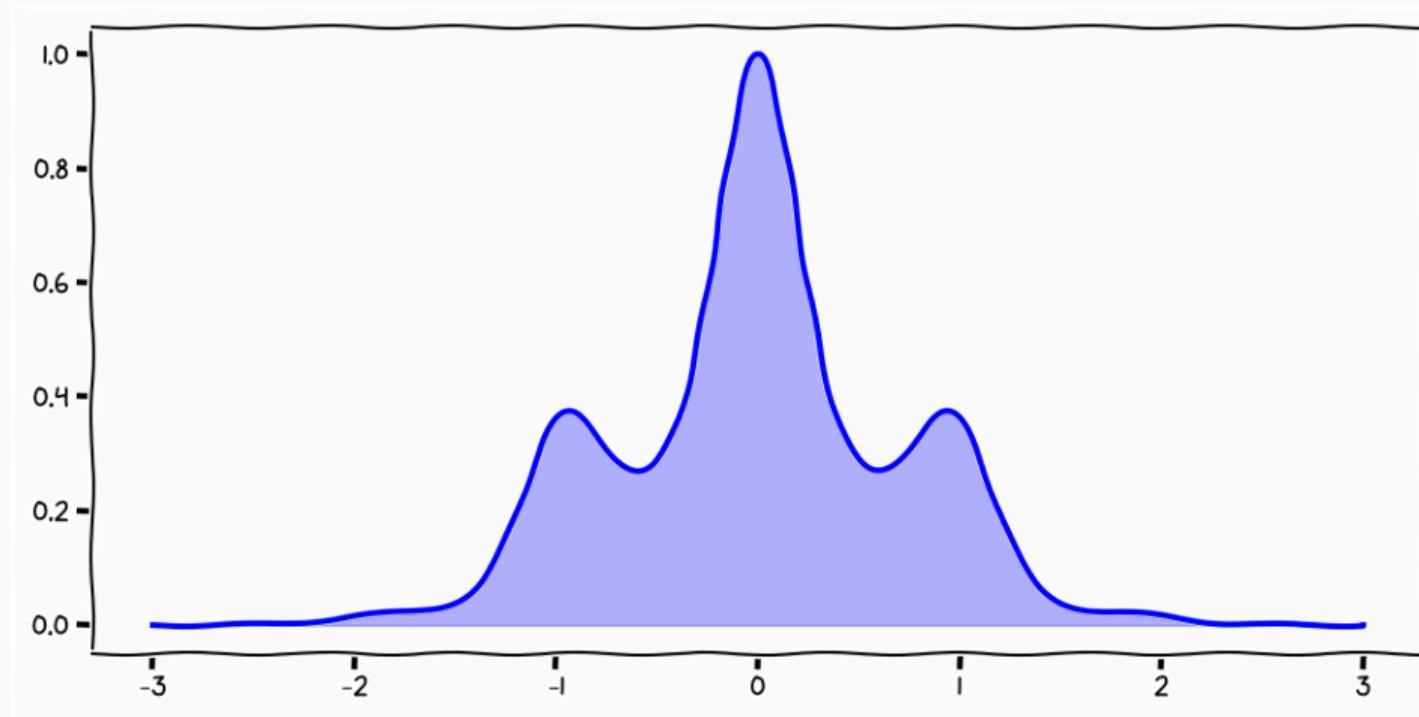
Bayesian Quadrature

Quantity of interest

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

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

Integration



What we would like

$$p(F \mid Y)$$

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

What we would like

$$p(F \mid Y)$$

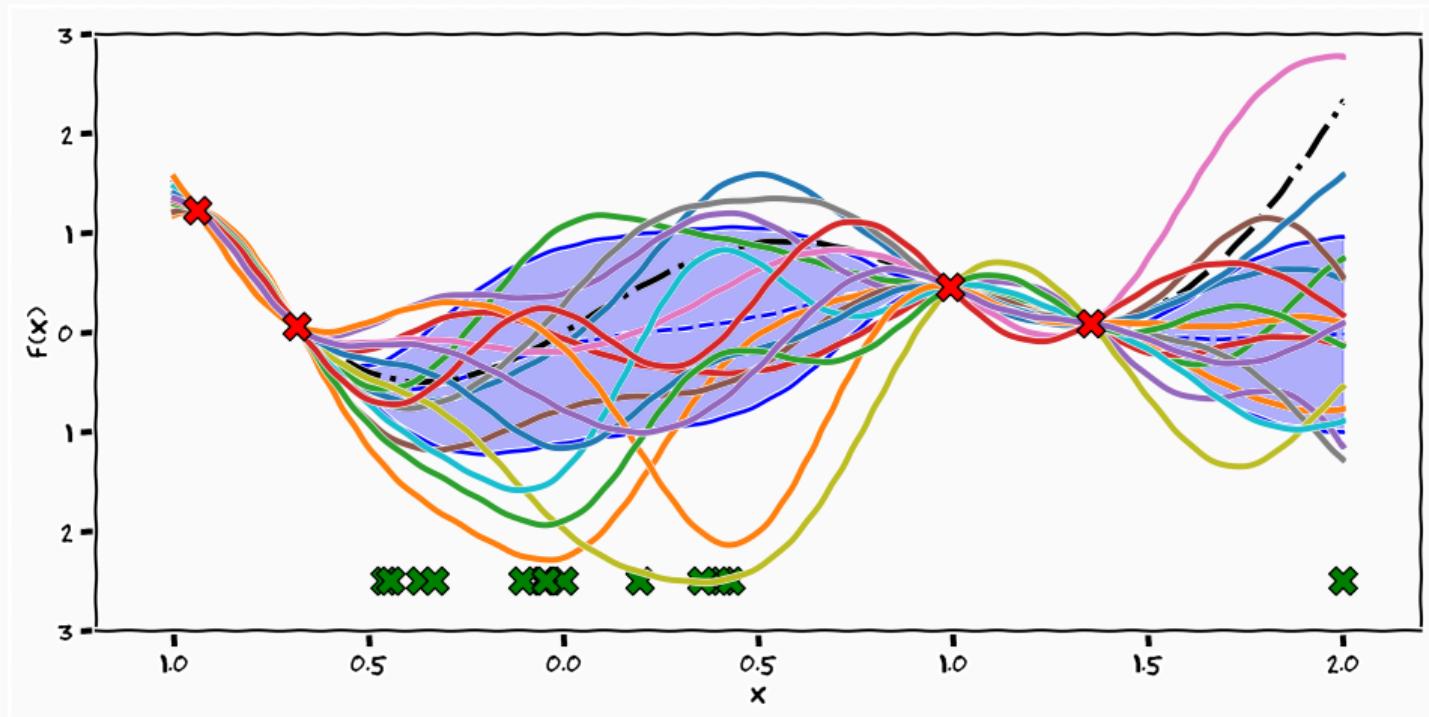
- given that I have seen data Y what is my belief about the integral
- allows for "active learning"

What we would like

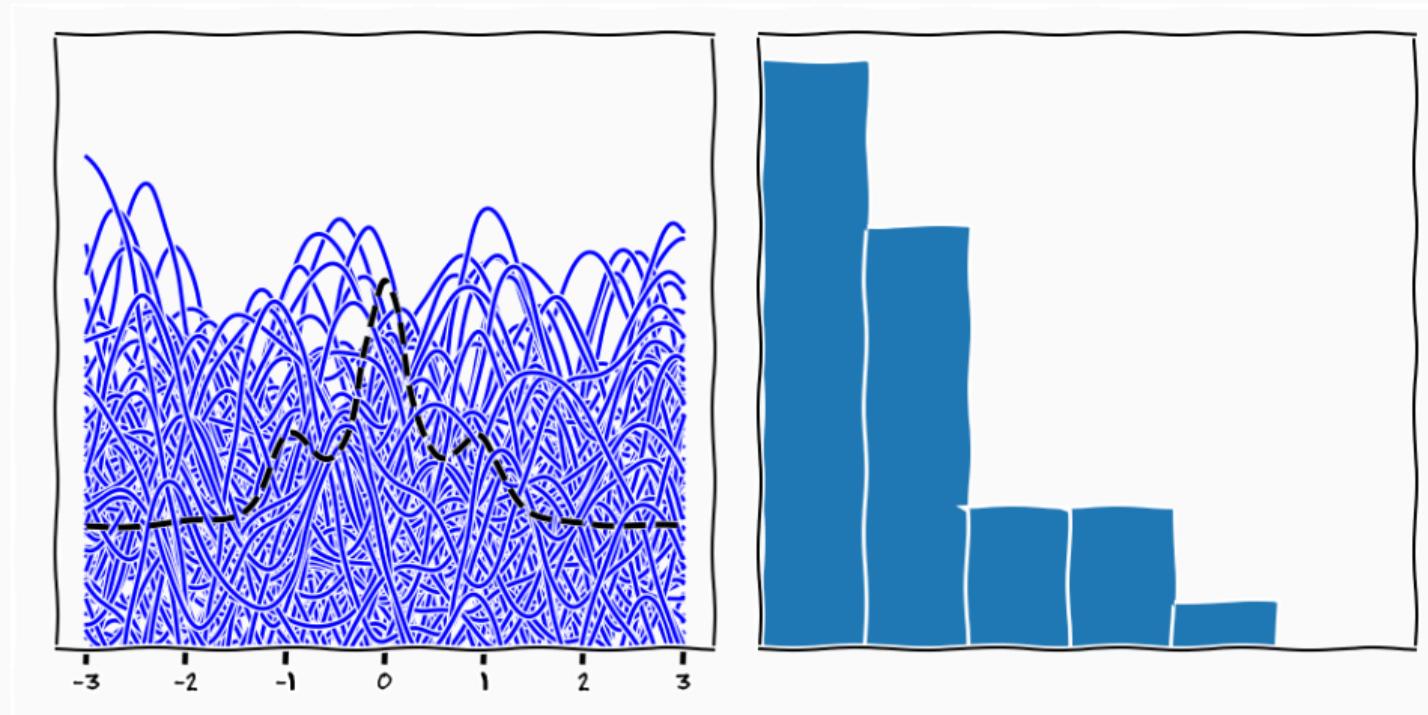
$$p(F \mid Y)$$

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

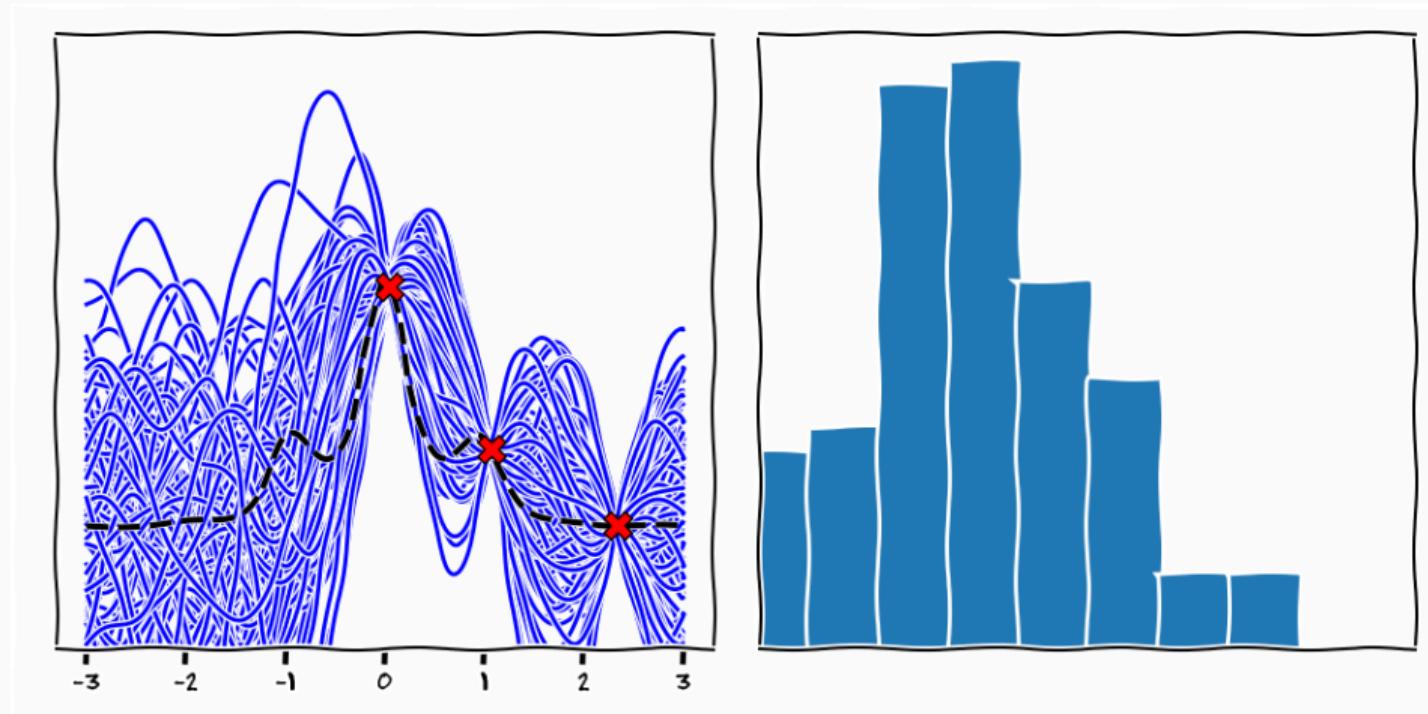
Emulation



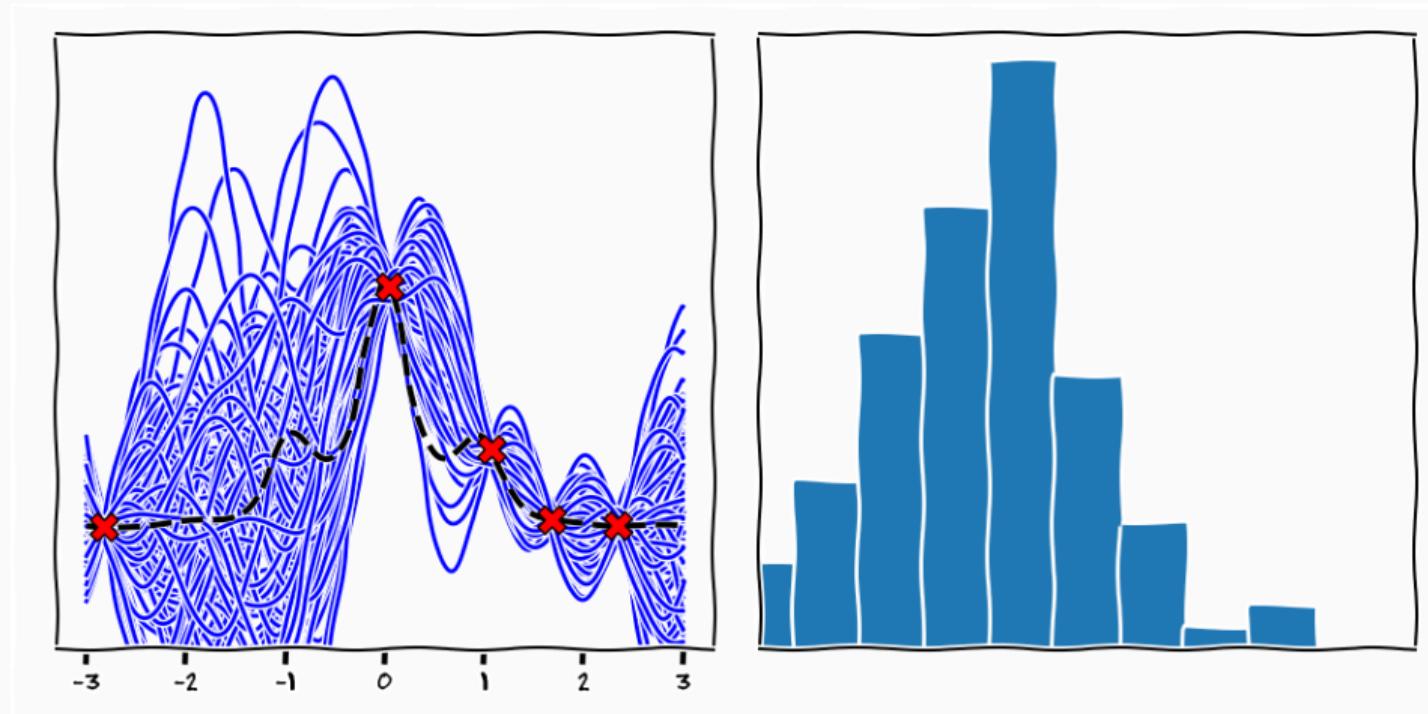
Quadrature



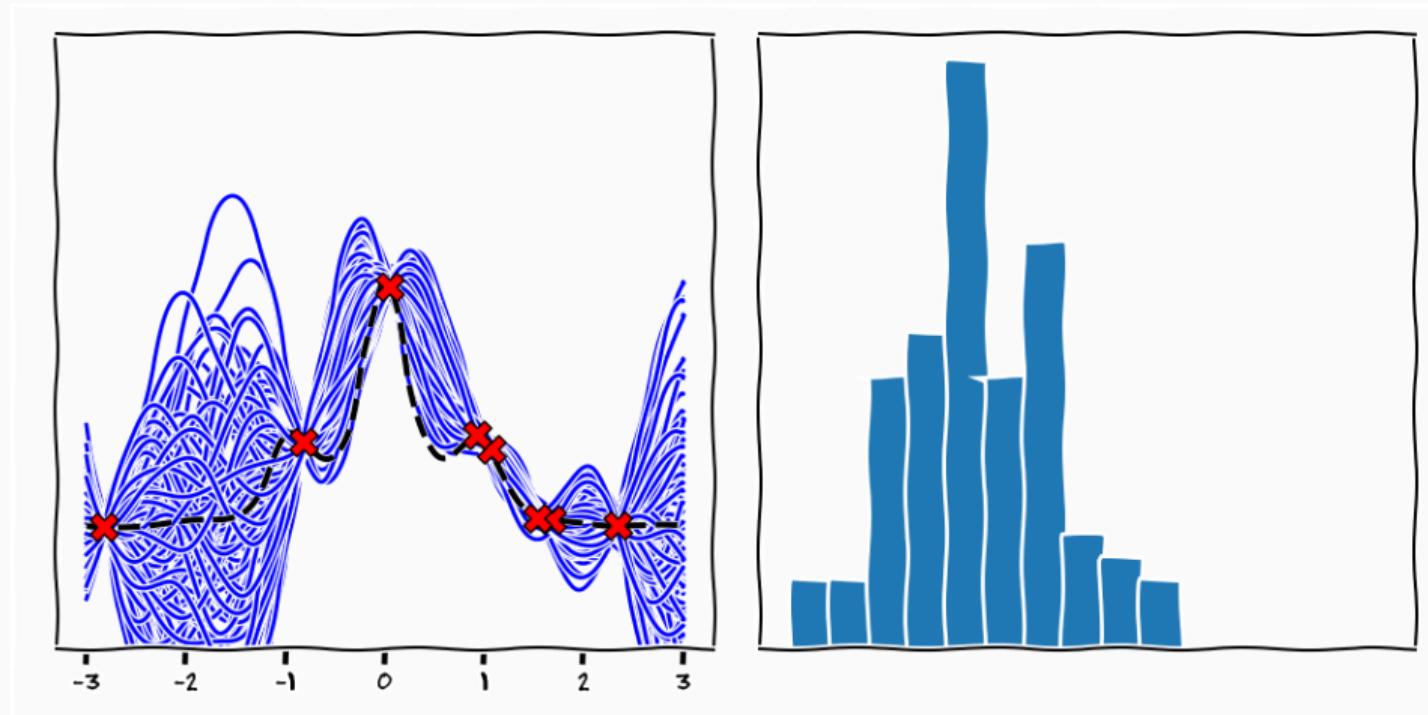
Quadrature



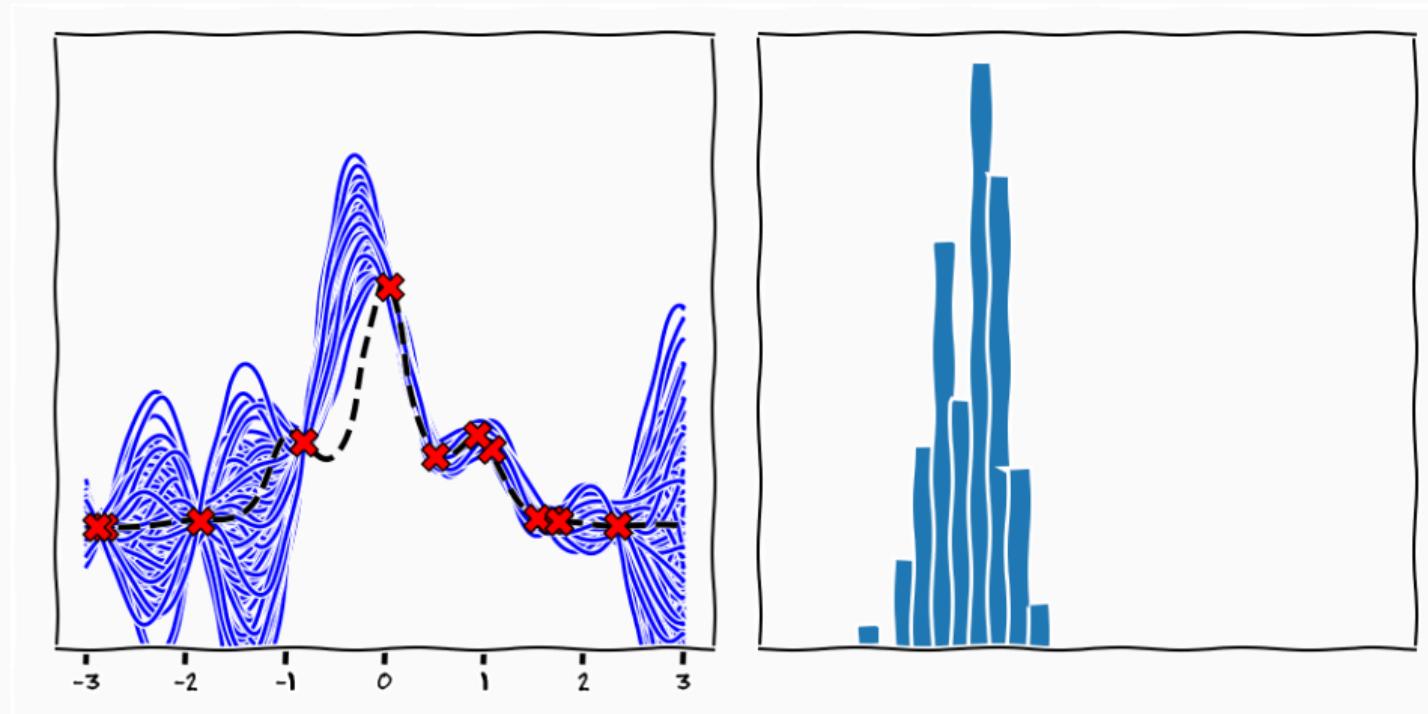
Quadrature



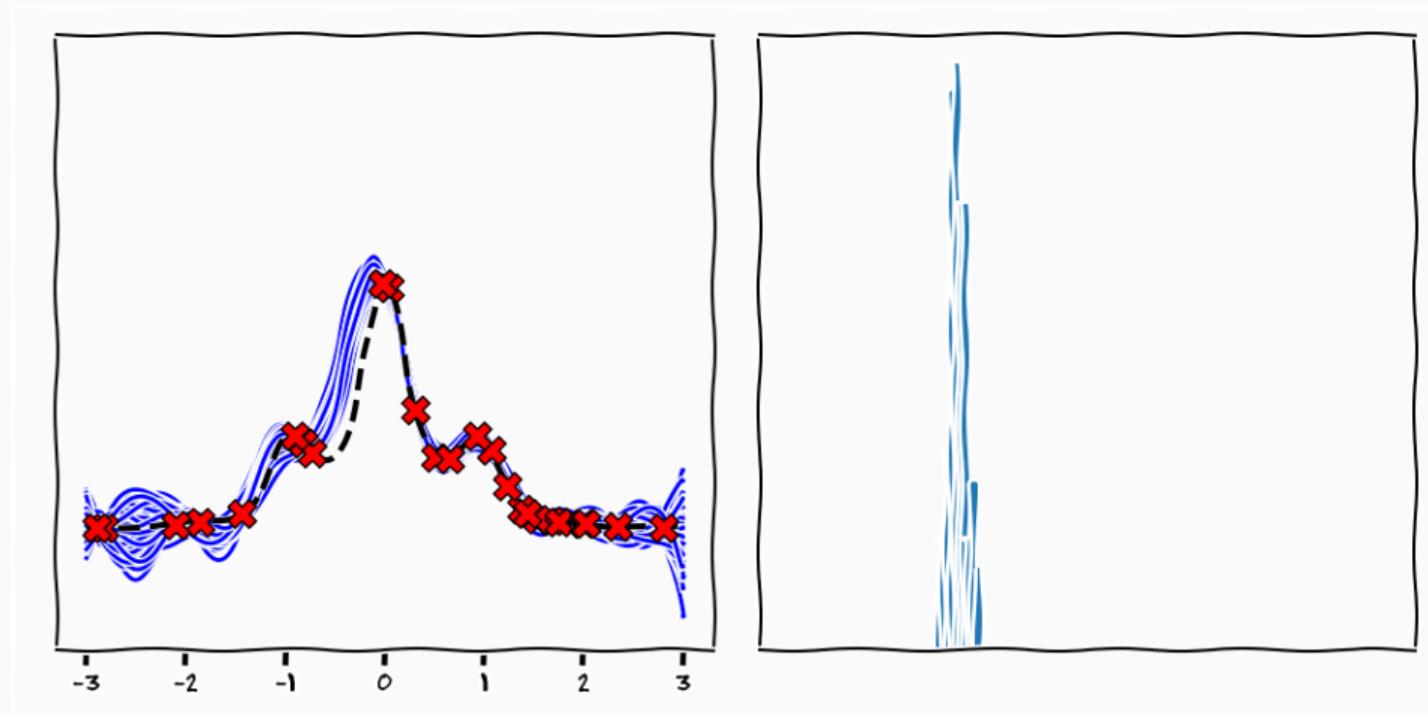
Quadrature



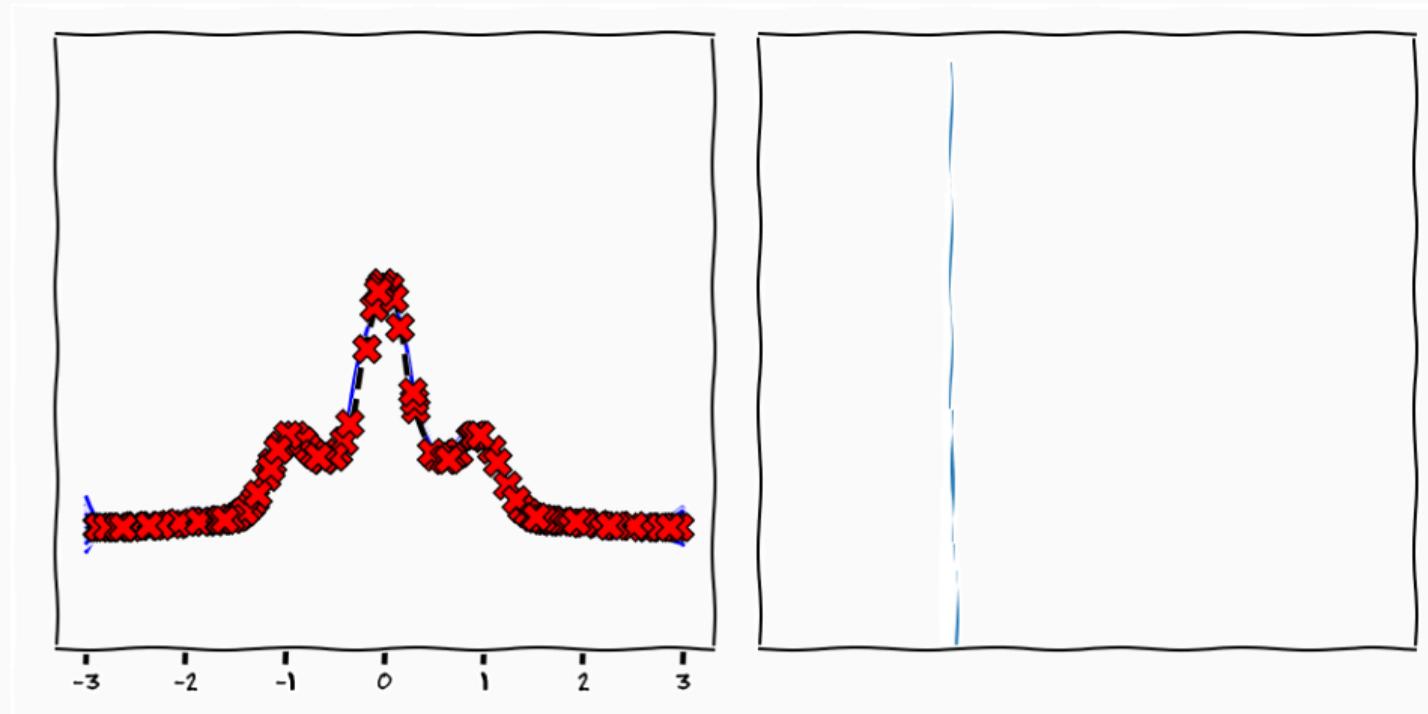
Quadrature



Quadrature



Quadrature



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

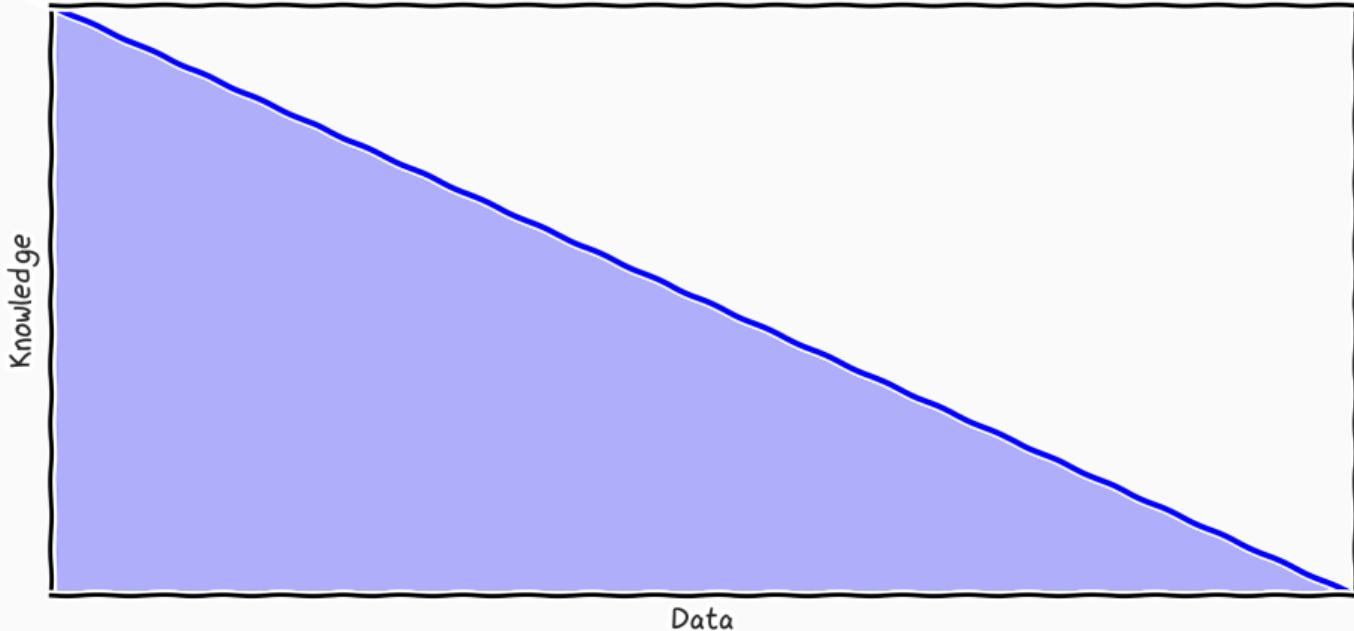
Knowledge

- $f(x)$ 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

Bayesian Quadrature [O'Hagan, 1991]

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

Bayesian Quadrature [O'Hagan, 1991]

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

Bayesian Quadrature [O'Hagan, 1991]

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

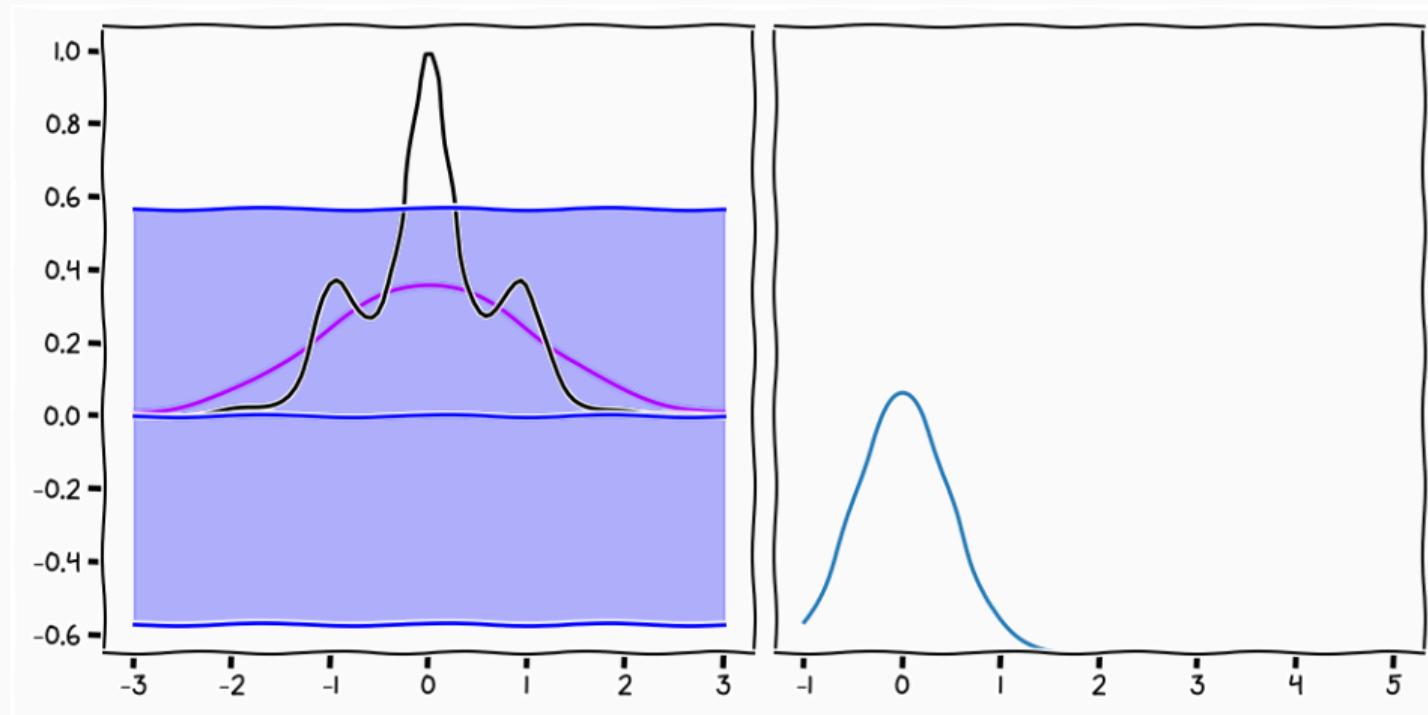
- We can derive $p(F | Y)$ through our normal conditioning procedure

Bayesian Quadrature [O'Hagan, 1991]

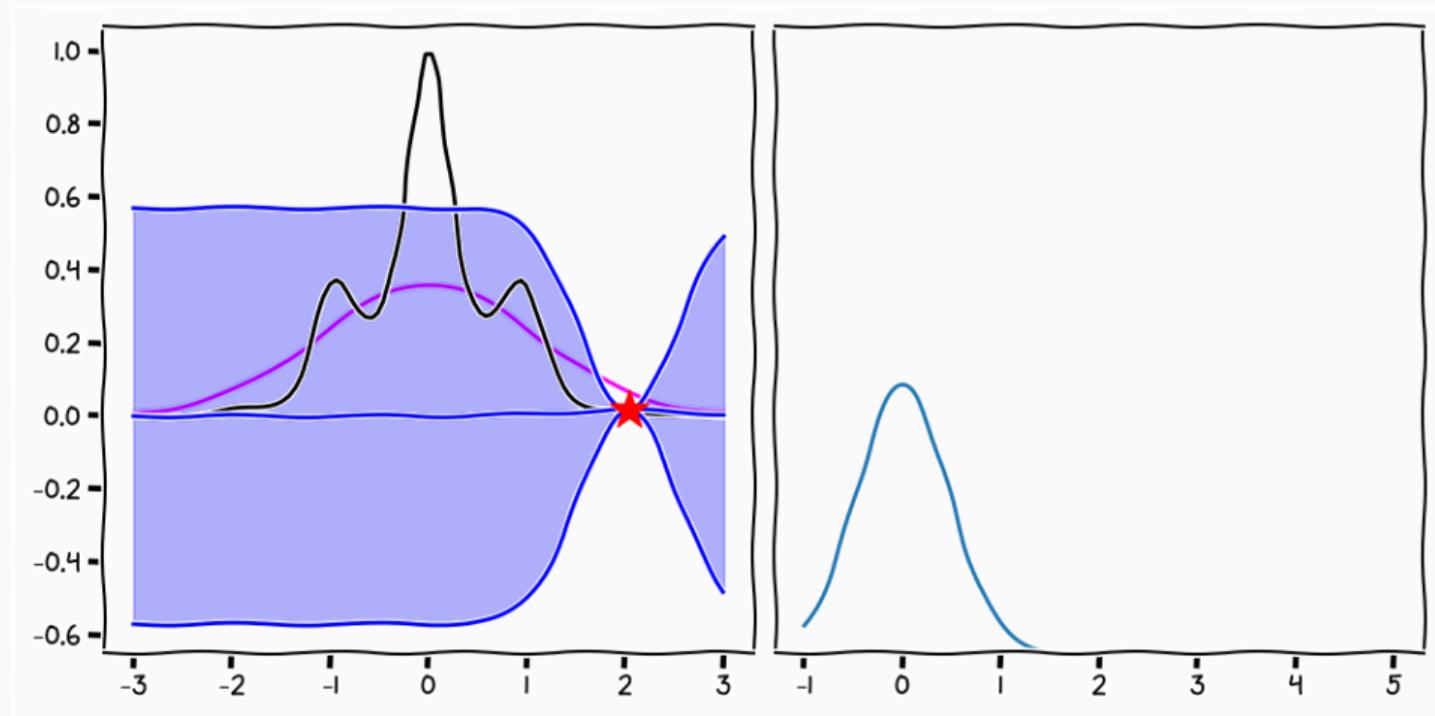
$$p \begin{pmatrix} Y \\ F \end{pmatrix} = \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_X \\ \int m_X(x)dx \end{bmatrix}, \begin{bmatrix} k(X, X) & \int k(X, x)dx \\ \int k(x, X)dx & \int \int k(x, x')dx dx' \end{bmatrix} \right)$$

- 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

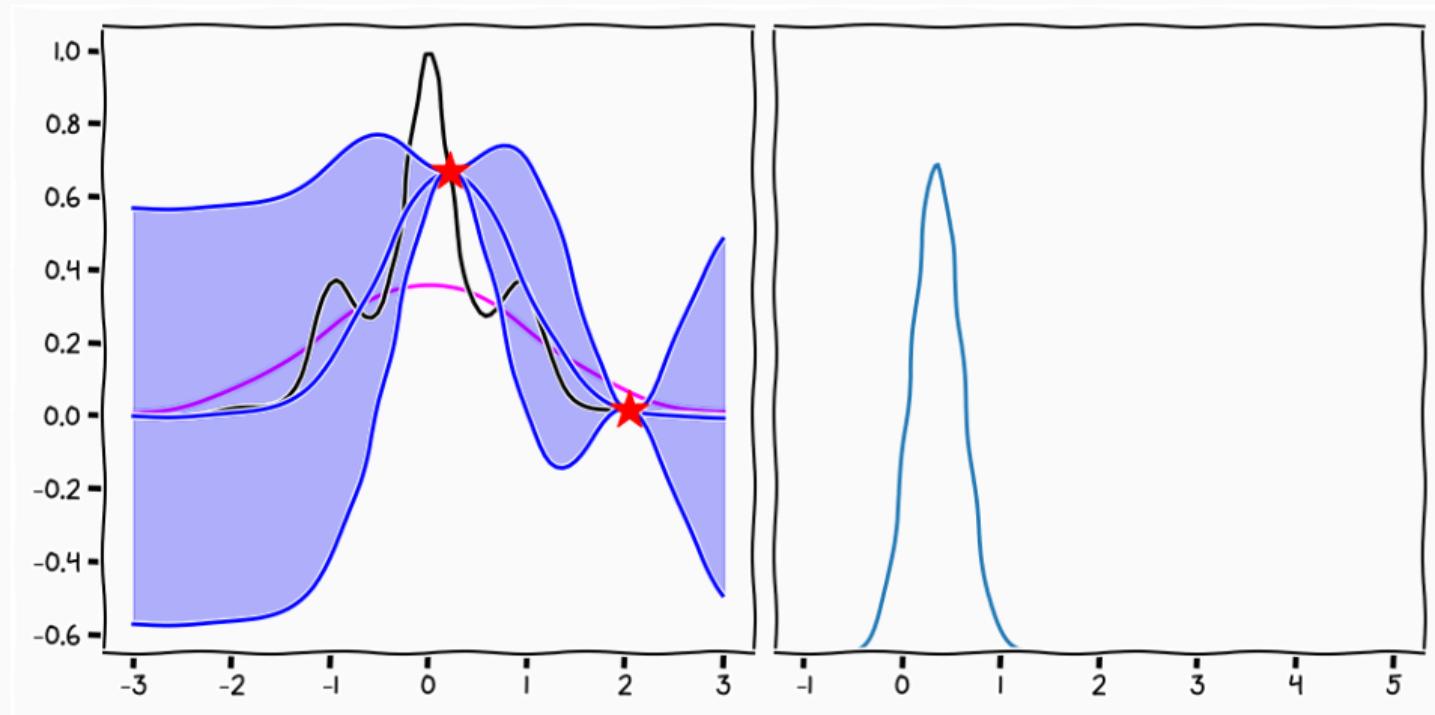
Statistical Inference



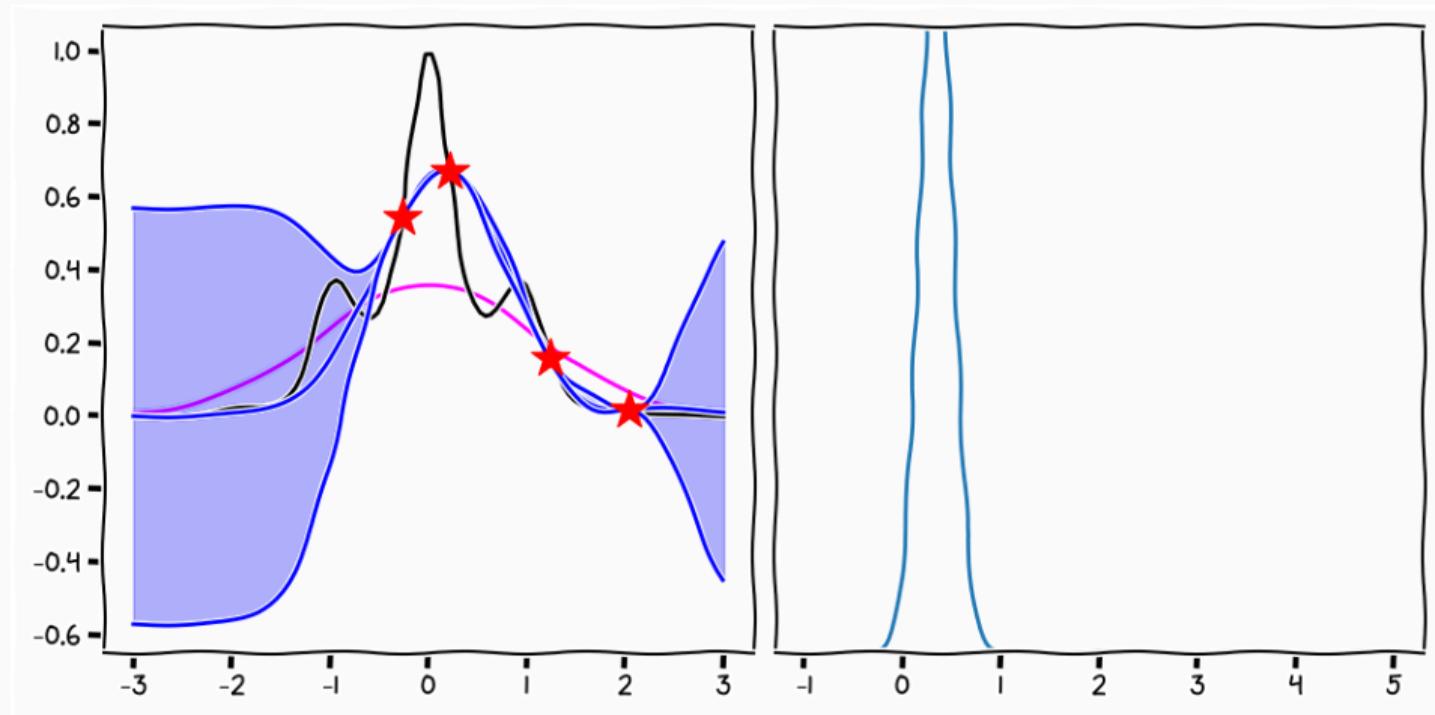
Statistical Inference



Statistical Inference



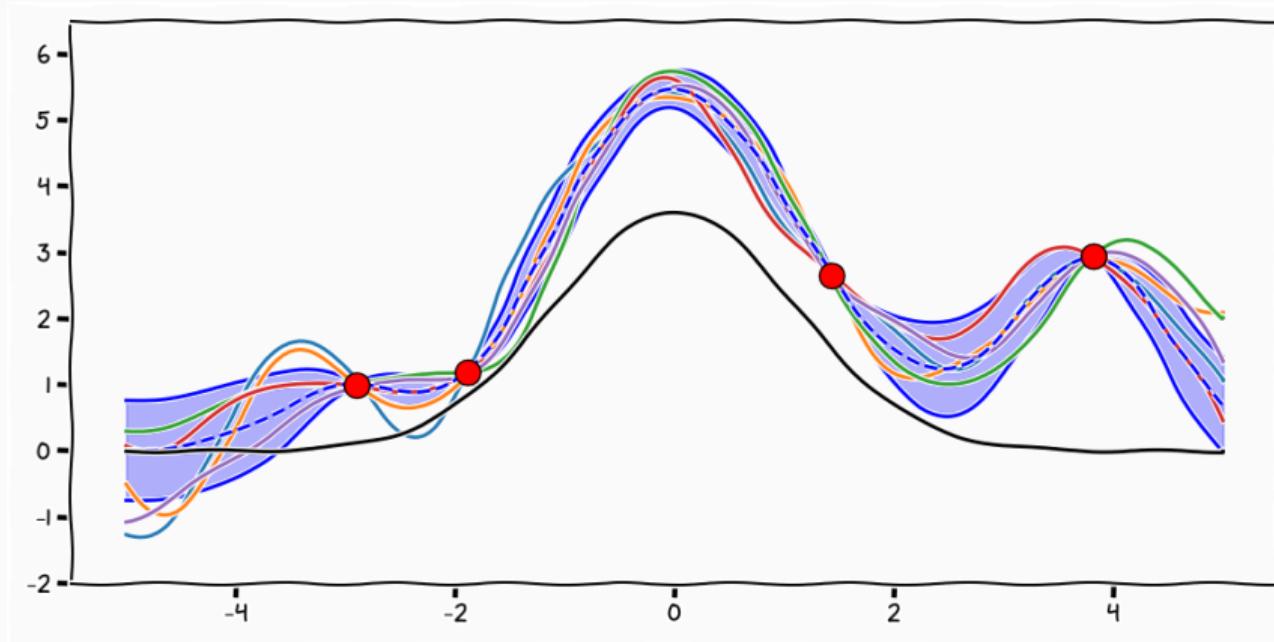
Statistical Inference



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

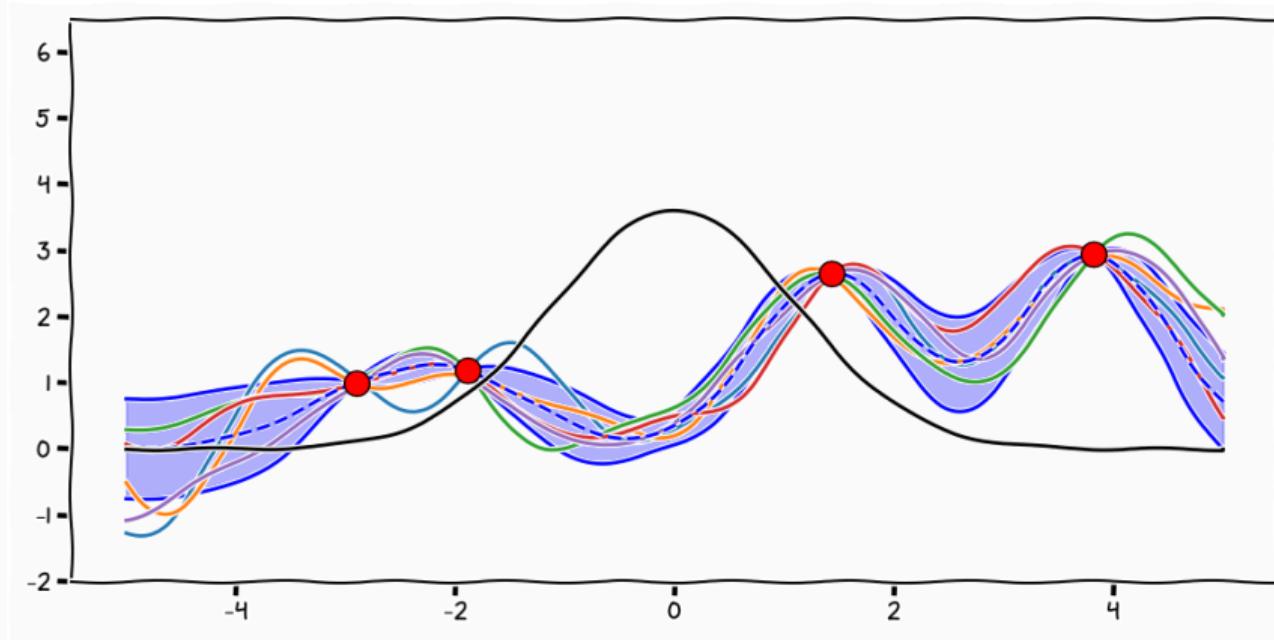
- We can derive $p(F | Y)$ through our normal conditioning procedure
- $p(F | Y) = \mathcal{N}(\mu_F, k_F)$ is a uni-variate Gaussian
- $p(Y | F) = \mathcal{N}(\mu_Y, k_Y)$ is a **Gaussian process**

Integral Constrained Samples



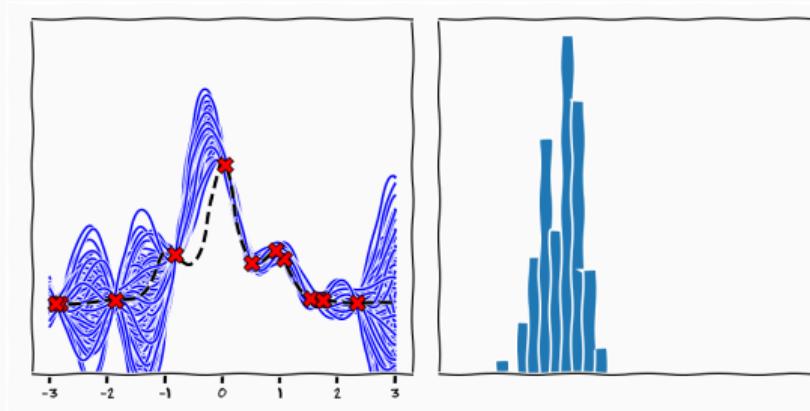
$$F = 4.0$$

Integral Constrained Samples



$$F = 1.0$$

Information Operator³

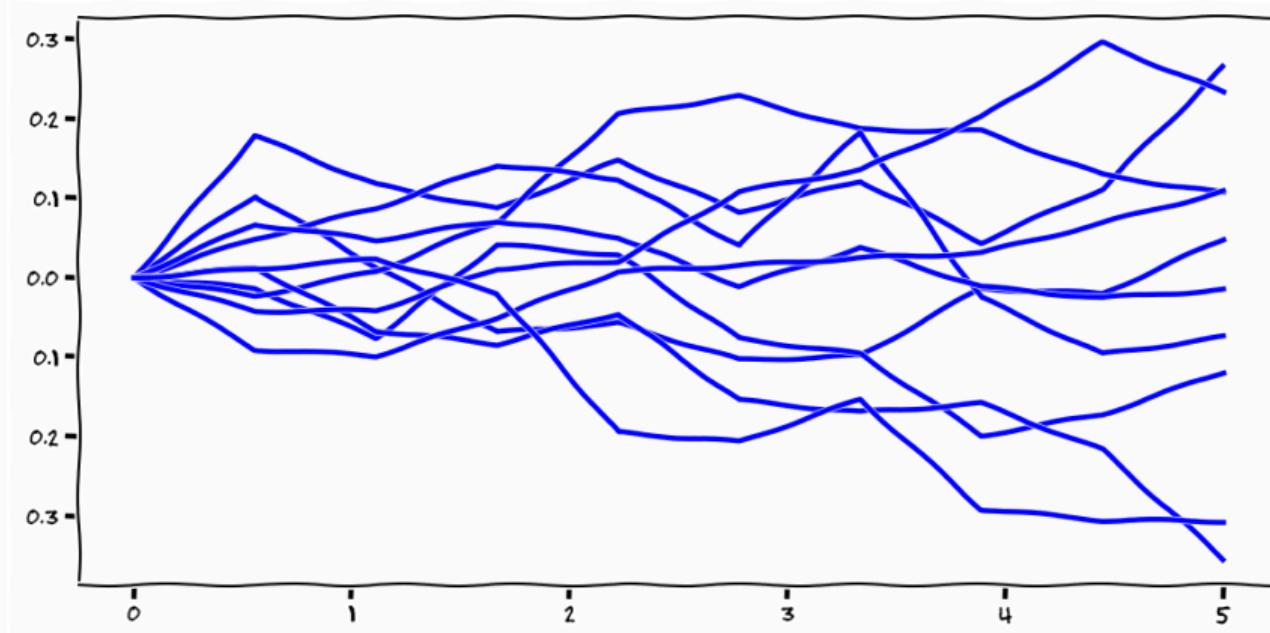


Integrand variance $\alpha(x) = k(x, x)$

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

³sometimes called a "Design Rule"

Choice of Covariance



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

Quadrature Rule

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

Quadrature 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!!!

Quadrature 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!!!!!!!

Trapedzoid Rule

Definition (Trapedzoid Rule)

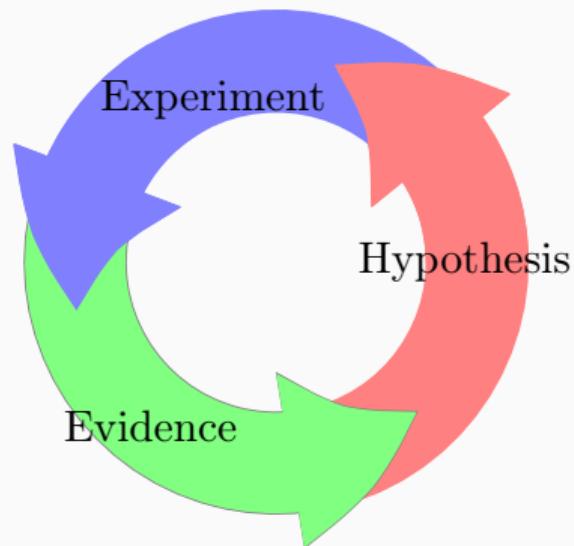
The trapezoidal rule is the posterior mean estimate for the integral

$F = \int_a^b f(x)dx$ 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

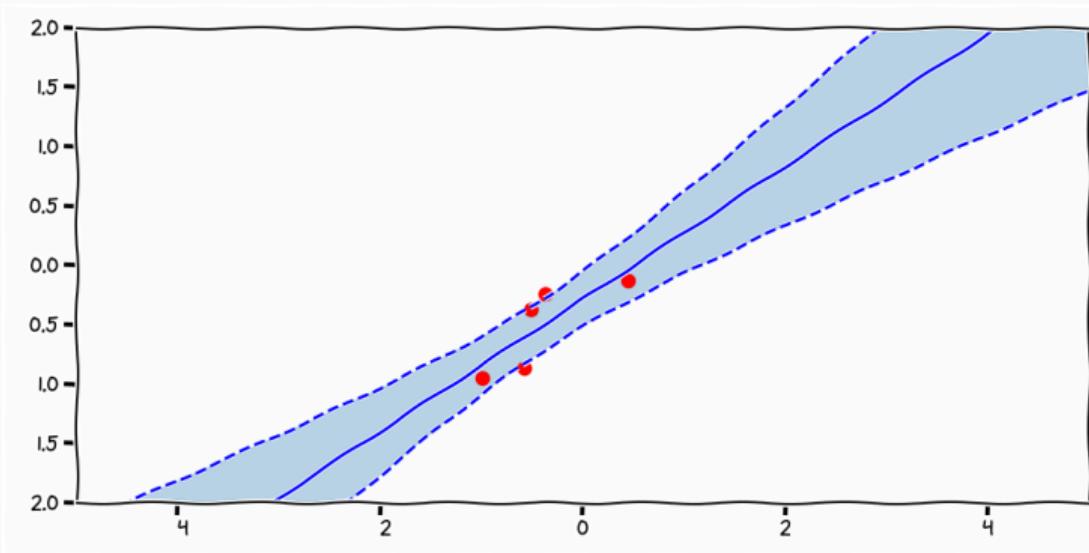


The Scientific Principle



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

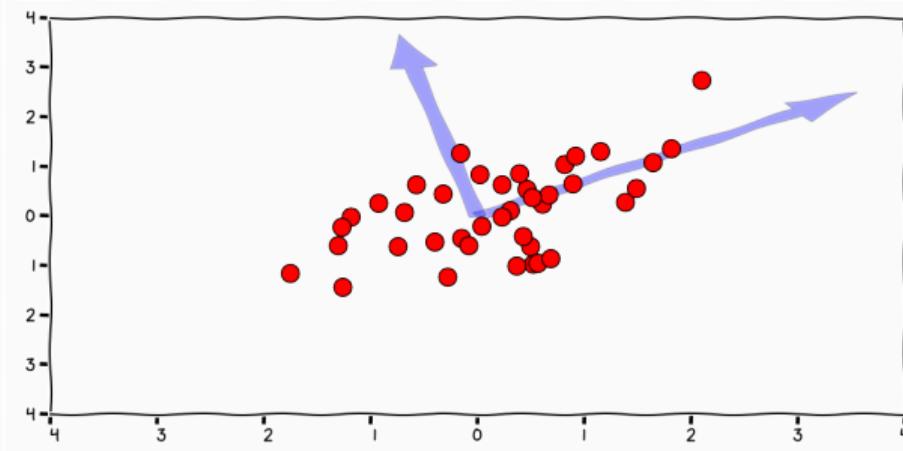
Least Squares Regression



Legendre (1805) algorithm that reduces "error"

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

Factor Analysis



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](#)

Why Probabilistic Numerics? `scipy.optimize.minimize`

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

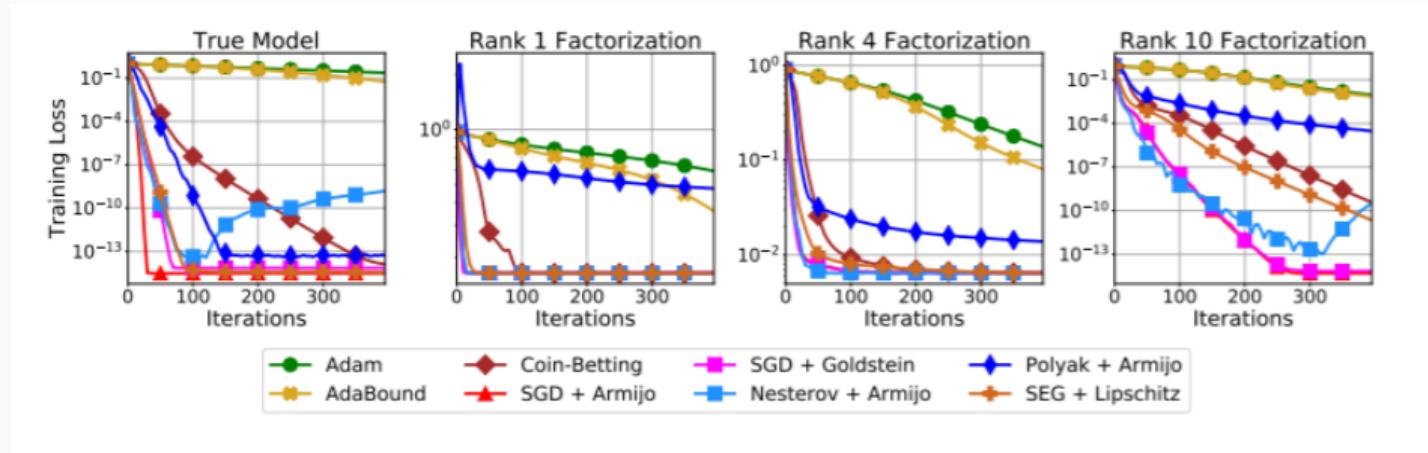
- There are tons of numerical algorithms for every problem under the sun
- They work really well
- They give different results on the same problem

- 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?*

Neural Networks (Maybe) Evolved to Make Adam The Best Optimizer

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.

No Free Lunch⁴



⁴Vaswani et al., 2019.

"The Machine Learning Principle"⁵

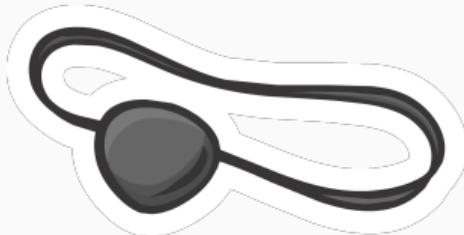
"There is a notion of success . . . which I think is novel in the history of science. It interprets success as approximating unanalyzed data."

– Prof. Noam Chomsky

⁵Chomsky et al., 1980



Assumptions: Algorithms



Statistical Learning

$$\mathcal{A}_{\mathcal{H}}(\mathcal{S})$$

Summary

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶

⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- Computation is the process of extracting a latent property, machine learning is the statistical process of updating beliefs about latent properties

⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- 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

⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- 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?

⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- 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

⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- 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

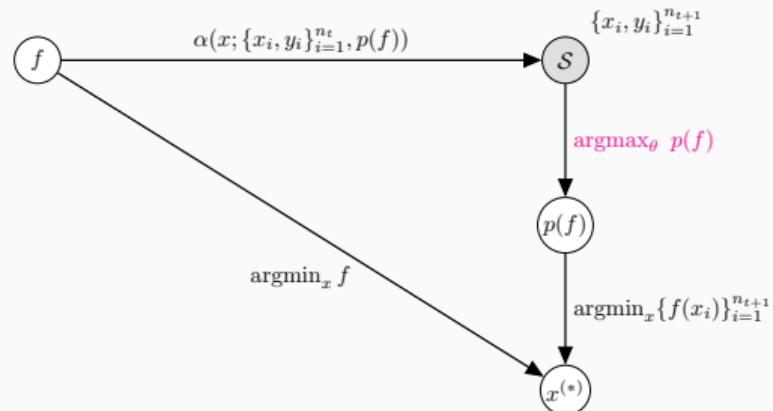
⁶these thoughts have been around for a long time

Summary

- Probabilistic Numerics extends the notion of statistical inference to computation⁶
- 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

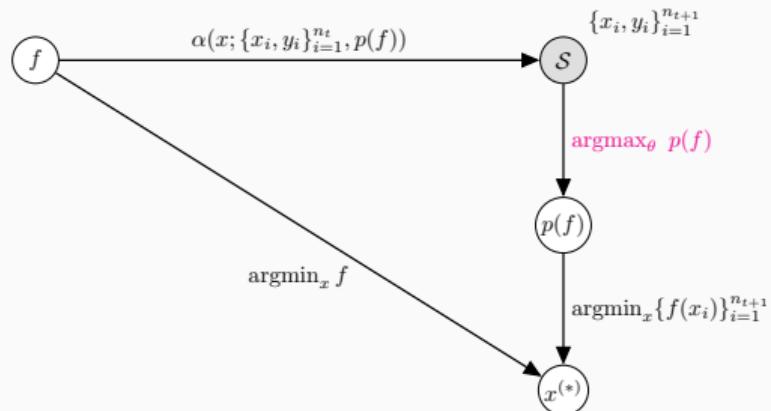
⁶these thoughts have been around for a long time

Is BO PN?



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

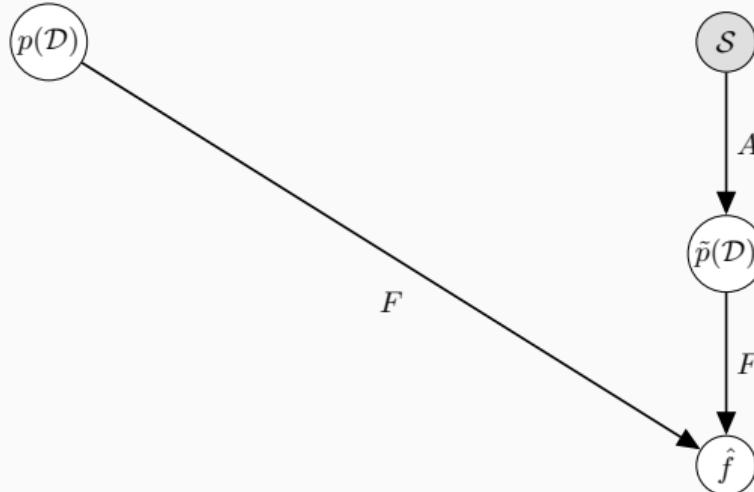
Is BO PN?



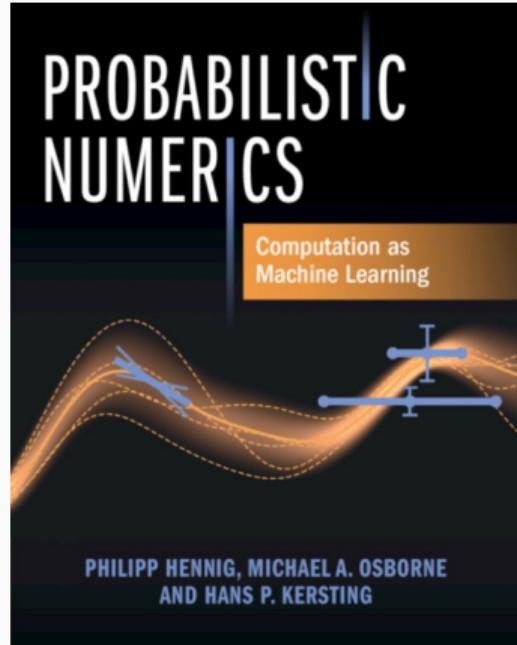
Yes it uses a probabilistic model as a proxy for decision loop

No the probabilistic model is not over the quantity of interest

Formalisation



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



<http://probnumschool.org>

eof

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.
- Lawrence, Neil D (2005). "Probabilistic non-linear principal component analysis with Gaussian process latent variable models." In: *Journal of Machine Learning Research* 6, pp. 1783–1816.
- Mackay, David J C (Dec. 1991). "Bayesian methods for adaptive models." PhD thesis. California Institute of Technology: California Institute of Technology.
- 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.
- ❑ Vaswani, Sharan et al. (2019). "Painless Stochastic Gradient: Interpolation, Line-Search, and Convergence Rates." In: *CoRR*.

Learning Hyper-parameters

$$\{\hat{\beta}, \hat{\ell}, \hat{\sigma}\} = \operatorname{argmax}_{\beta, \ell, \sigma} \log \int p(y | f) p(f) df$$

- Update the **hyper-parameters** of the GP
- We can do this by gradient based optimisation
 - *or with a Bayes Opt loop!*

⁷...and here I was thinking that you guys had some principles in life

Learning Hyper-parameters

$$\{\hat{\beta}, \hat{\ell}, \hat{\sigma}\} = \operatorname{argmax}_{\beta, \ell, \sigma} \log \int p(y | f) p(f) df$$

- Update the **hyper-parameters** of the GP
- We can do this by gradient based optimisation
 - *or with a Bayes Opt loop!*
- Are we doing **MLE**⁷?

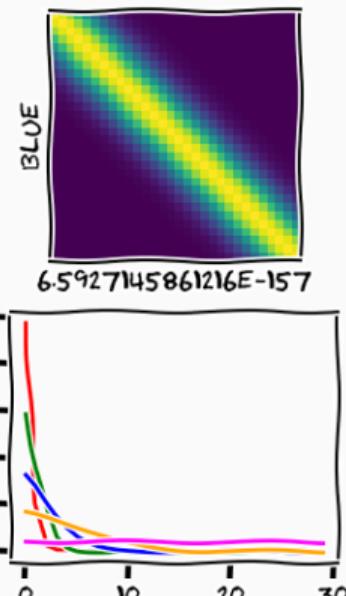
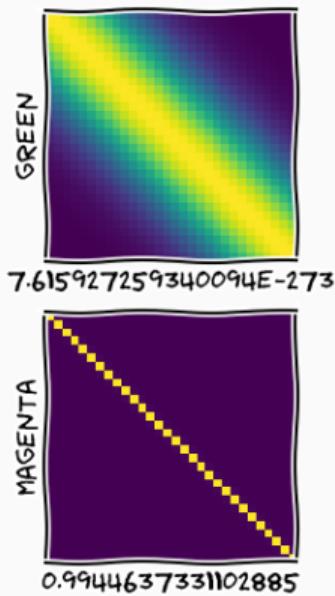
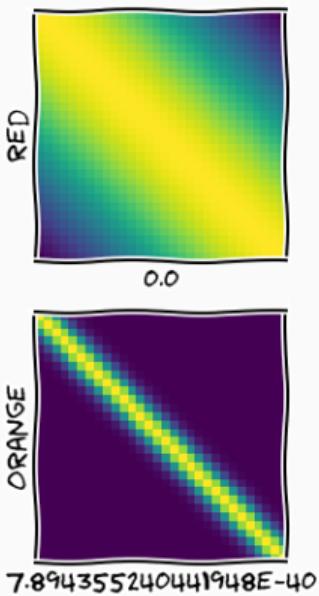
⁷...and here I was thinking that you guys had some principles in life

Marginal Likelihood

$$\begin{aligned}\{\hat{\beta}, \hat{\ell}, \hat{\sigma}\} &= \underset{\beta, \ell, \sigma}{\operatorname{argmax}} \log \int p(y \mid f) p(f) d f \\&= \underset{\beta, \ell, \sigma}{\operatorname{argmin}} -\log p(y) \\&= \underset{\beta, \ell, \sigma}{\operatorname{argmin}} \frac{1}{2} \text{trace}(\mathbf{Y} \mathbf{K}^{-1} \mathbf{Y}^T) + \frac{1}{2} \log |K| + \frac{N}{2} \log(2\pi)\end{aligned}$$

- *Data – fit* how well does the observations fit the model
- *”Complexity”* how “smooth” is the functions

Determinant



How to implement



Semantics

$$p(\theta | y) = \frac{p(y | \theta)p(\theta)}{\int p(y | \theta)p(\theta)d\theta}$$

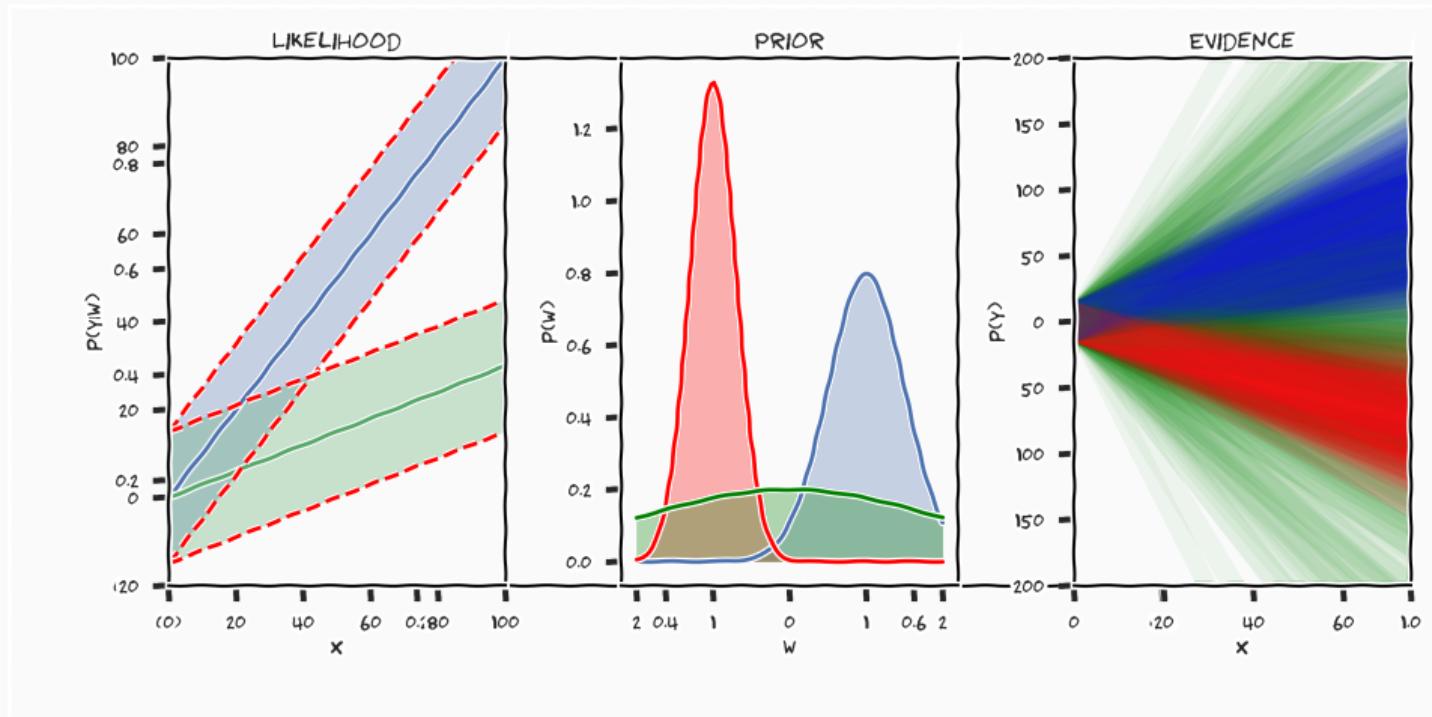
Likelihood How much **evidence** is there in the data for a specific hypothesis

Prior What are my beliefs about different hypothesis

Posterior What is my **updated** belief after having seen data

Evidence What is my belief about the data

Regression Model



Marginalisation



*Next time you want to give your friends a compliment, tell them that you have completely **marginalised** them from your life*

Regression Models

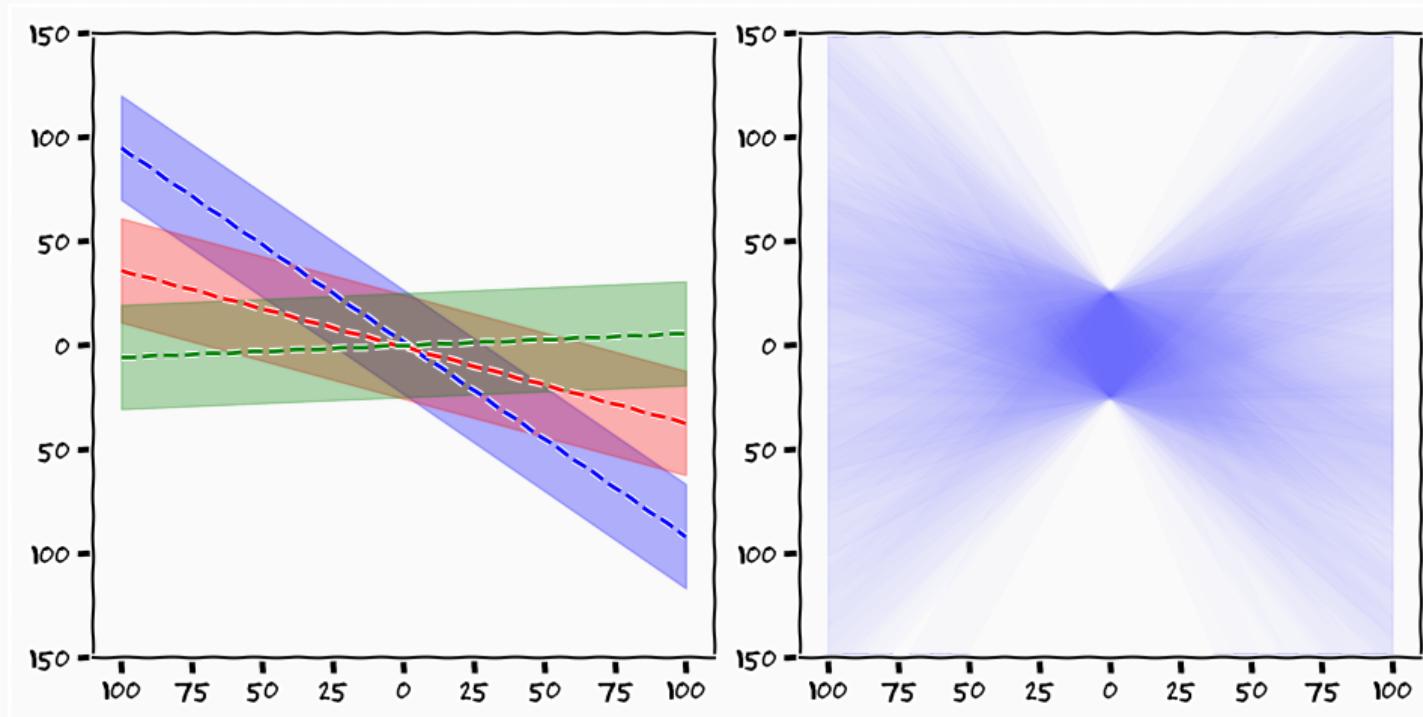
Linear Model

$$p(y_i|x_i, \mathbf{w}) = \mathcal{N}(w_0 + w_1 \cdot x_i, \beta^{-1})$$

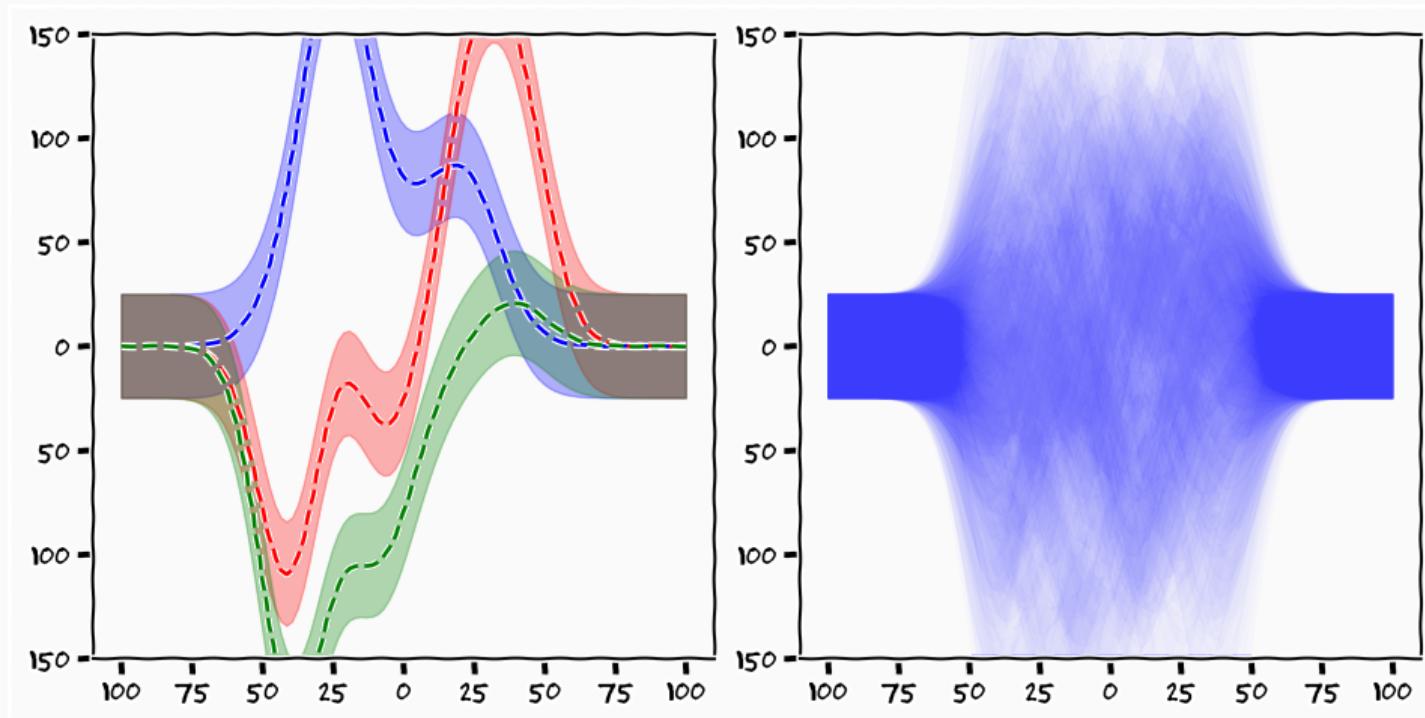
Basis function

$$p(y_i|x_i, \mathbf{w}) = \mathcal{N}\left(\sum_{i=1}^6 w_i \phi(x_i), \beta^{-1}\right)$$

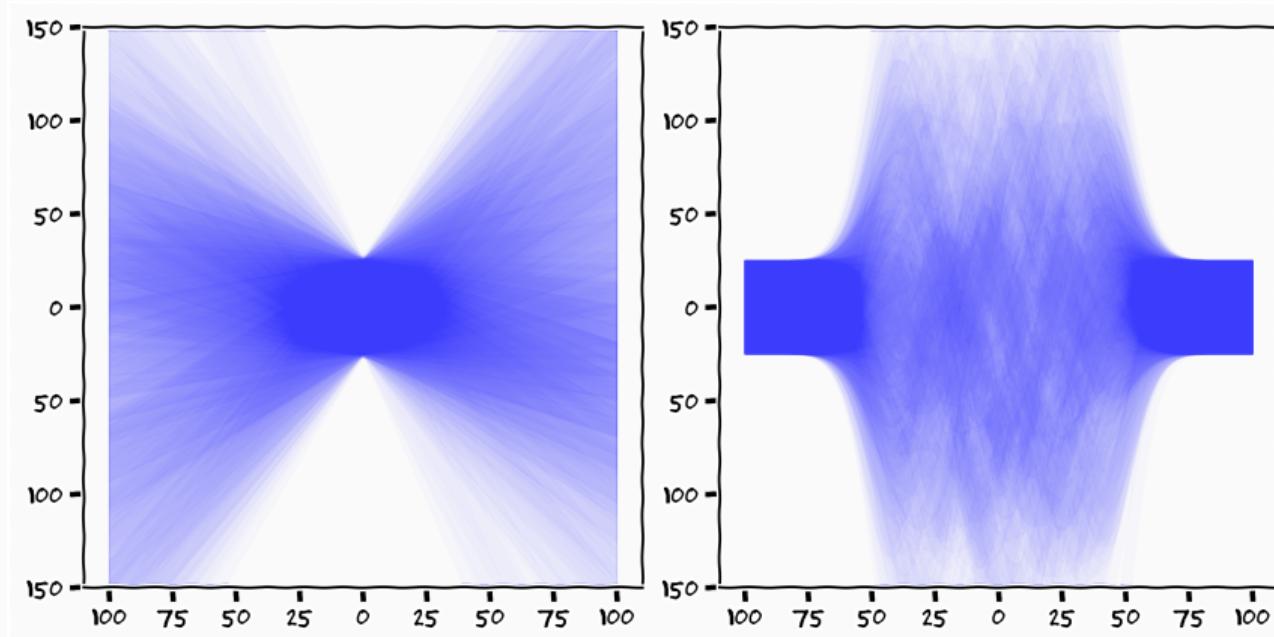
Linear Linear Regression



Linear Regression



Evidence

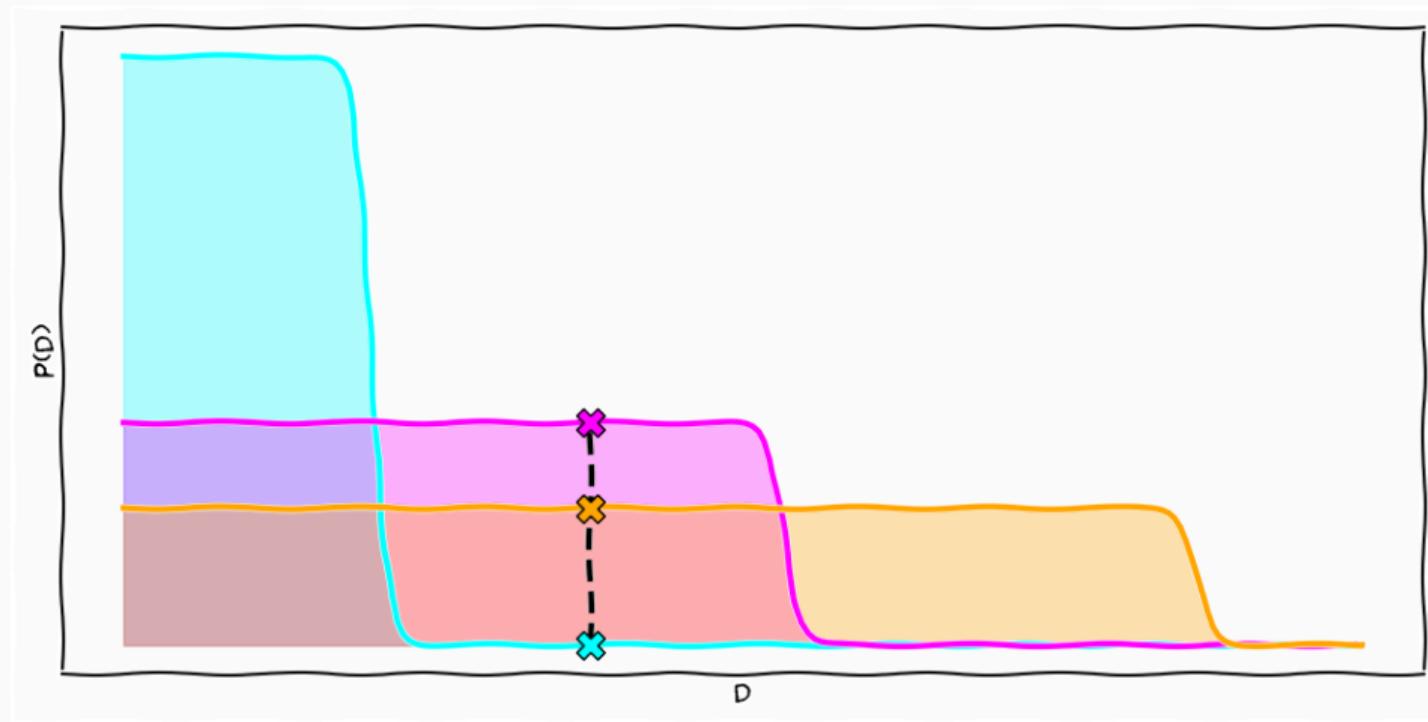


$$p(\mathcal{Y}) = \int p(\mathcal{Y}|\mathbf{W})p(\mathbf{W})d\mathbf{W}$$

Probabilities are a zero-sum game



The MacKay Plot Mackay, 1991



Occams Razor



Definition (Occams Razor)

"All things being equal, the simplest solution tends to be the best one"

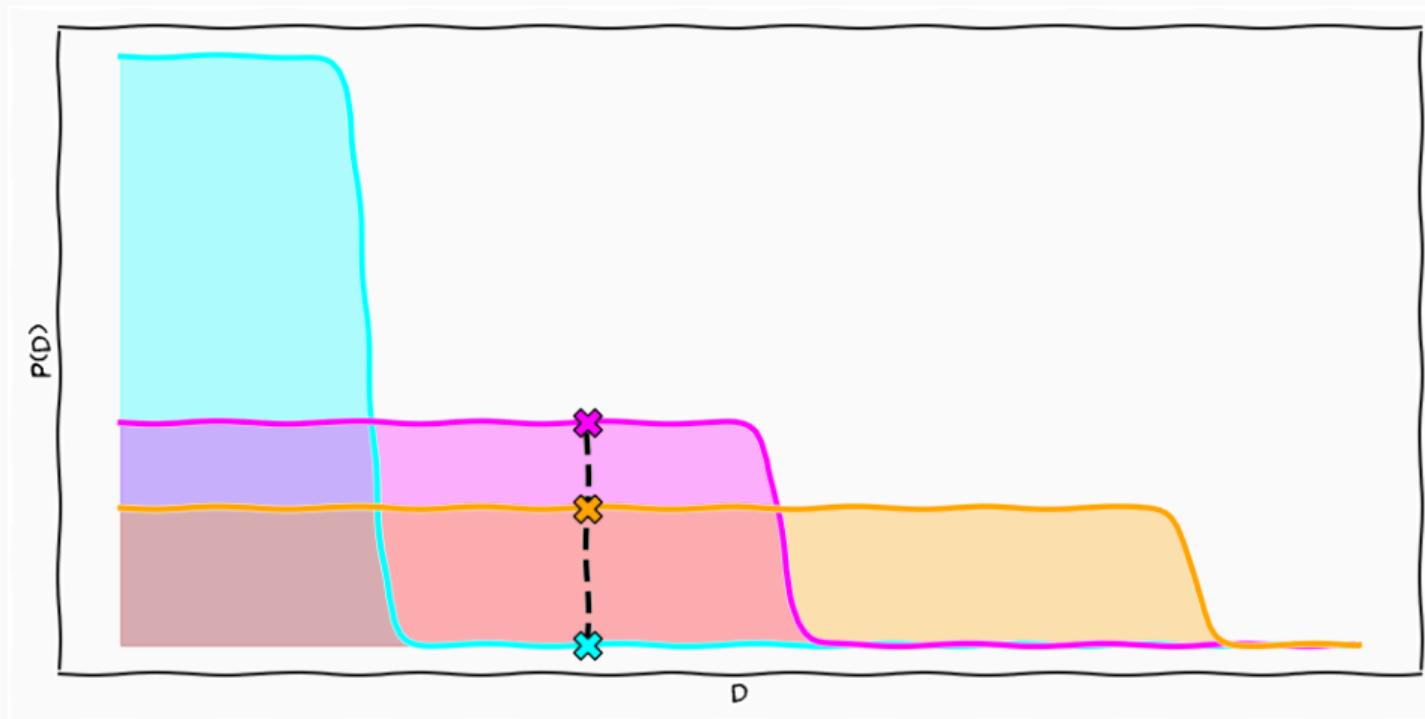
– William of Ockham

What is Simple?⁸



⁸<https://www.imdb.com/title/tt8132700/>

The MacKay Plot Mackay, 1991



Unsupervised Learning [Lawrence, 2005]

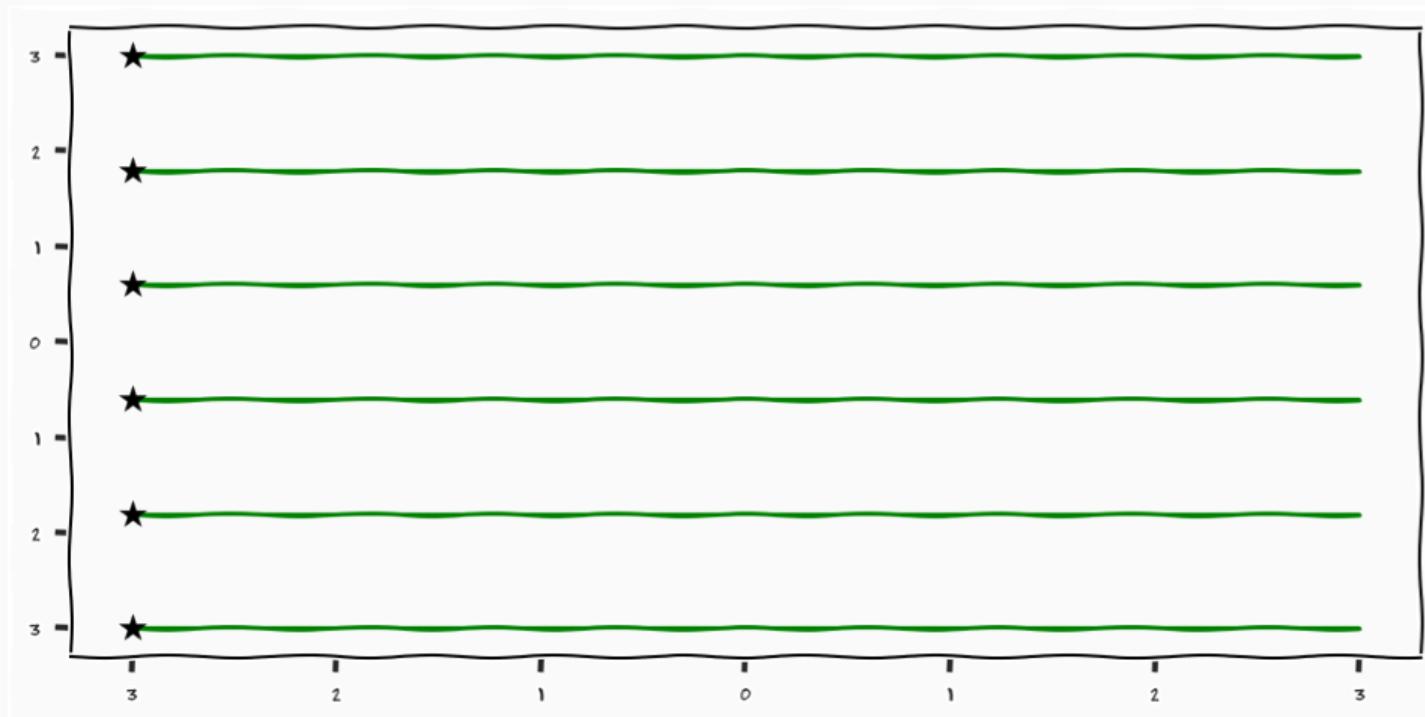
- Regression

$$p(y \mid x) = \int p(y \mid f)p(f \mid x)df$$

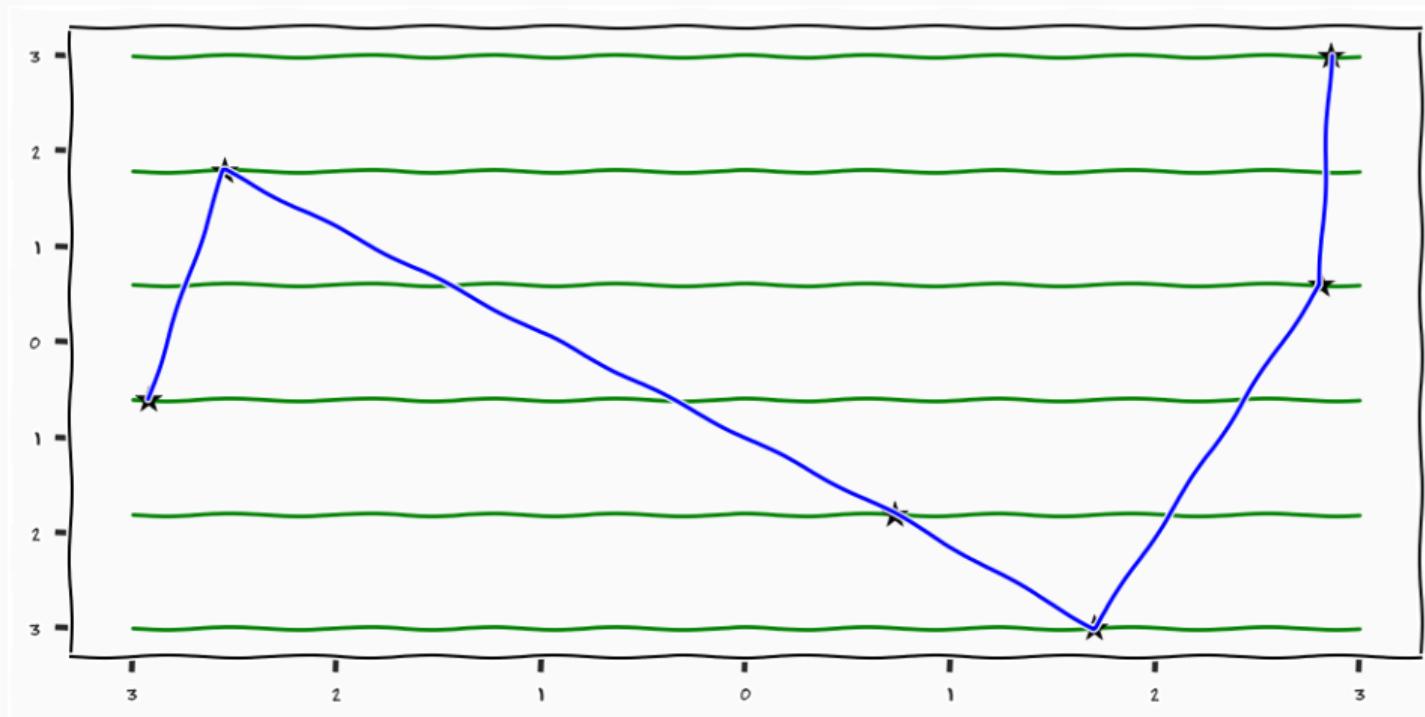
- "Unsupervised" Learning

$$p(y) = \int p(y \mid f)p(f \mid x)p(x)dfdx$$

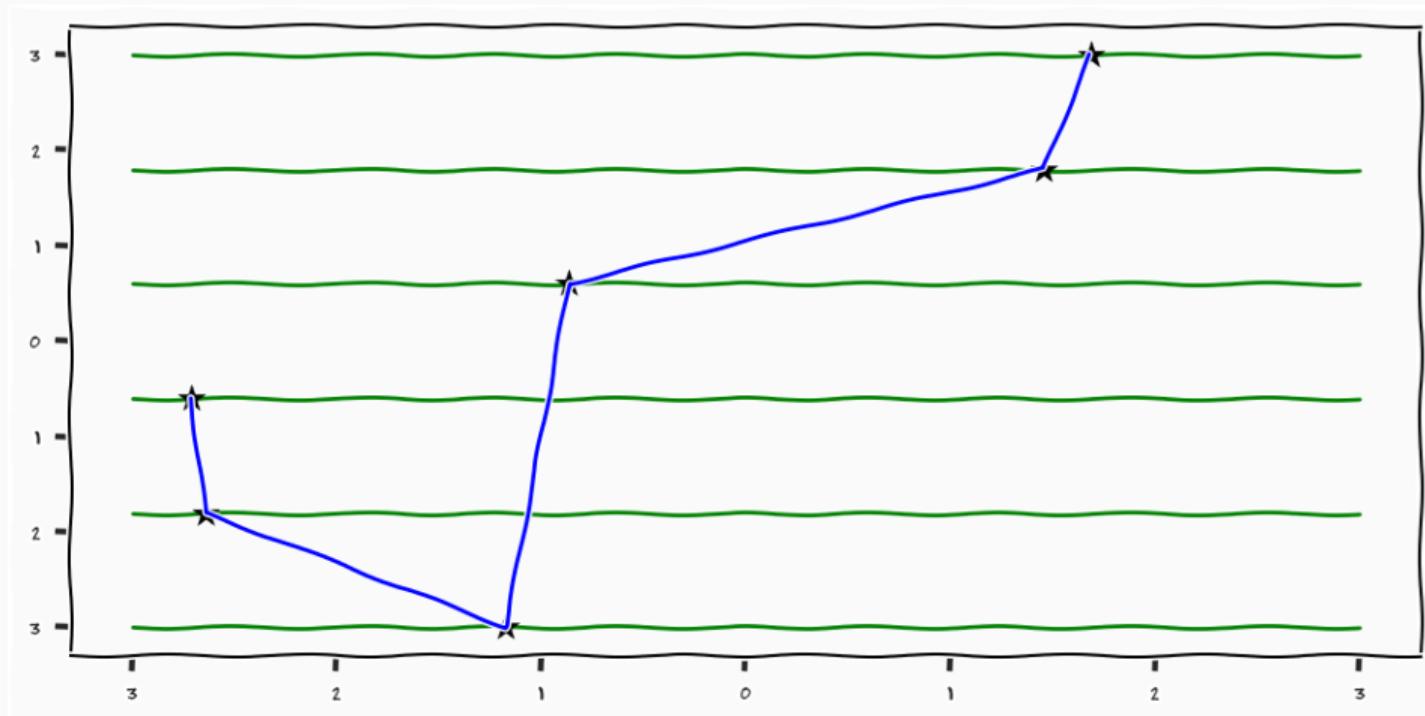
Unsupervised Learning



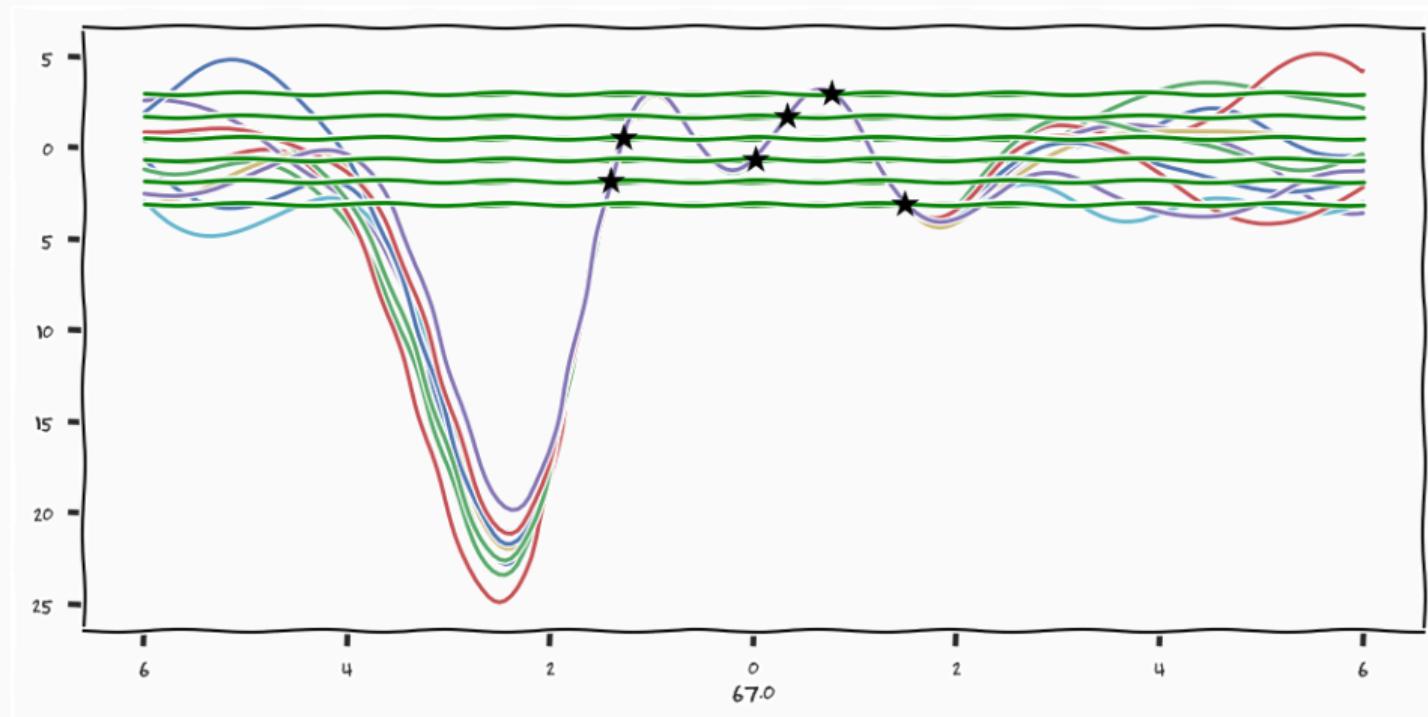
Unsupervised Learning



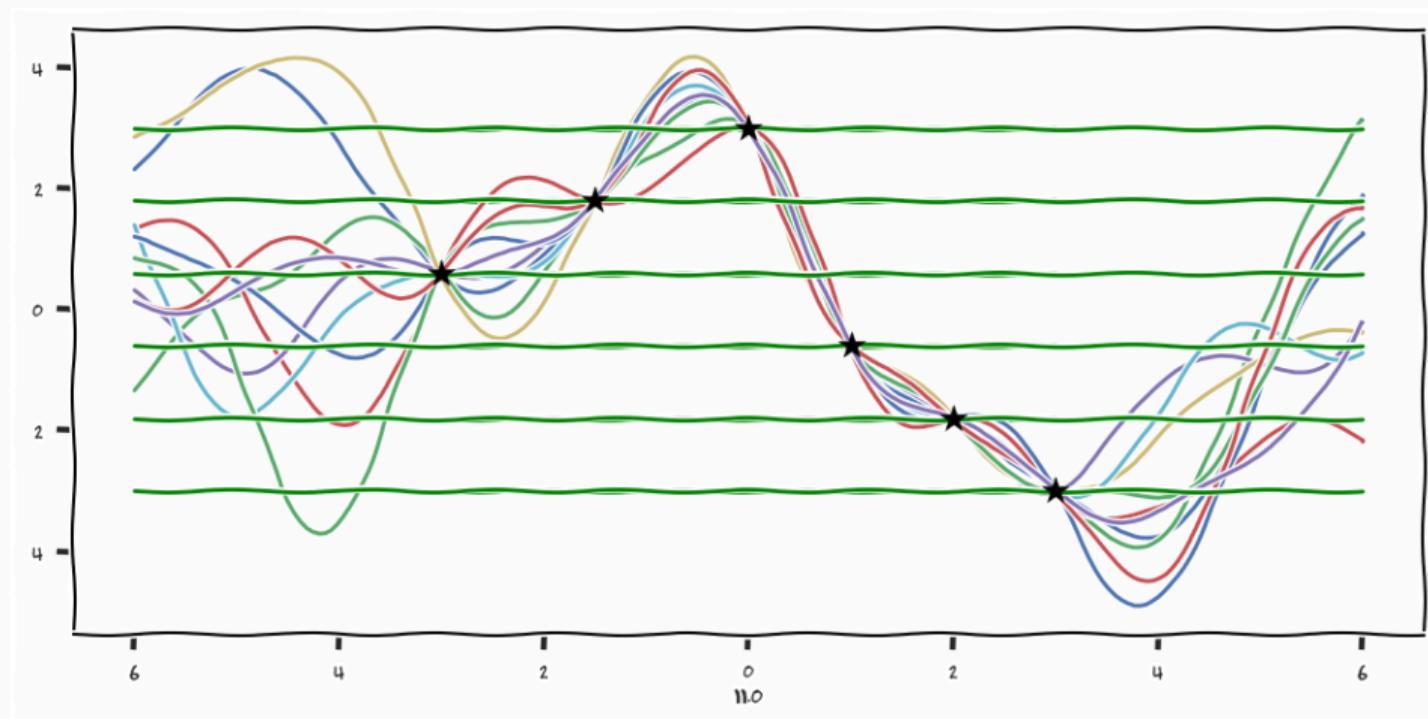
Unsupervised Learning



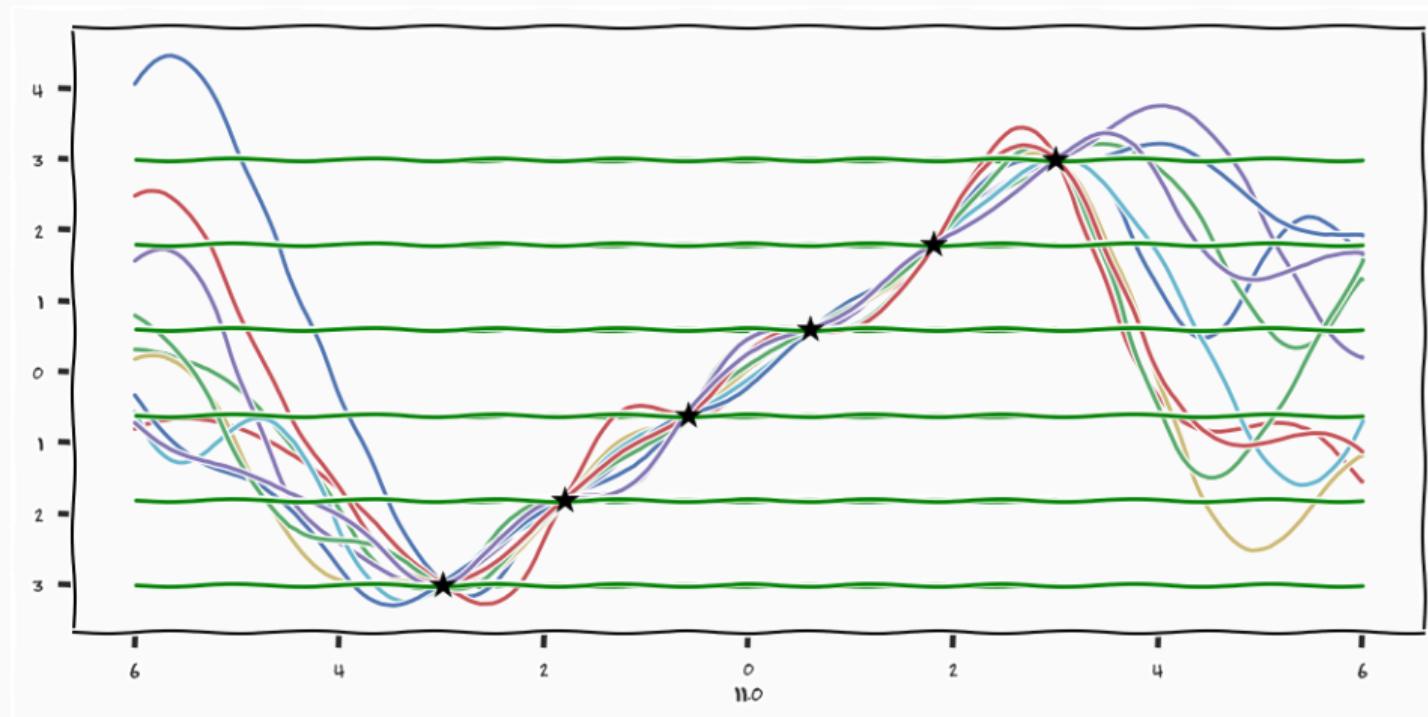
Gaussian Process Latent Variable Model [Lawrence, 2005]



Gaussian Process Latent Variable Model [Lawrence, 2005]



Gaussian Process Latent Variable Model [Lawrence, 2005]



Gaussian Process Latent Variable Model [Lawrence, 2005]

