

Probabilistic numerics for deep learning

Presenter: Shijia Wang

Michael Osborne

Department of Engineering Science, University of Oxford

Deep Learning (DLSS) and Reinforcement Learning (RLSS) Summer
School, Montreal 2017

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- Integration beats Optimization

3 Conclusion

- Experiments

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- Integration beats Optimization

3 Conclusion

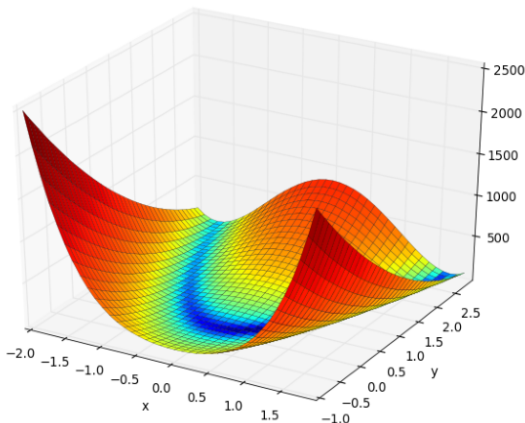
- Experiments

Definition

- Take the things we were most interested in achieving and apply to computation
- Apply probability theory to numerics (computation cores)

- Use numeric functions as learning algorithms
- Idea is to use Bayesian probability theories

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$



- Easy to graph on a computer
- No easy way of finding its global optimum
- Reason: computational limits from the optimization problem

- Epistemically uncertain about the function due to being unable to afford computation
- Probabilistically model function and use tools from decision theory to make optimal use of computation

Outline

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- Integration beats Optimization

3 Conclusion

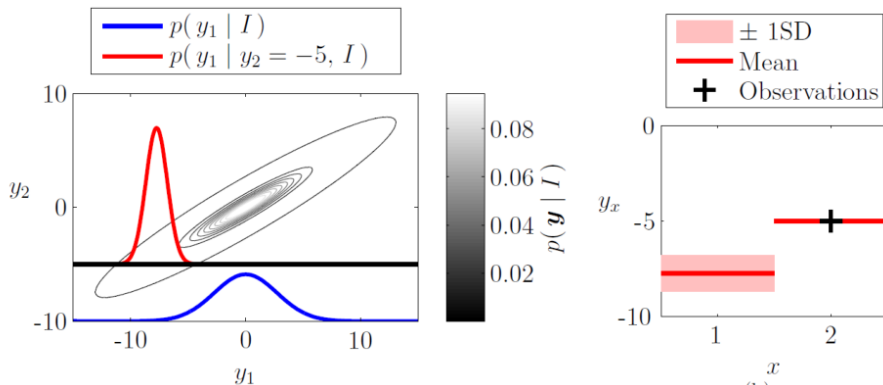
- Experiments

Probability

- Probability is an expression of confidence in a proposition
- Probability theory can quantify inverse of logic expression
- Depends on the agent's prior knowledge

Gaussian Distribution

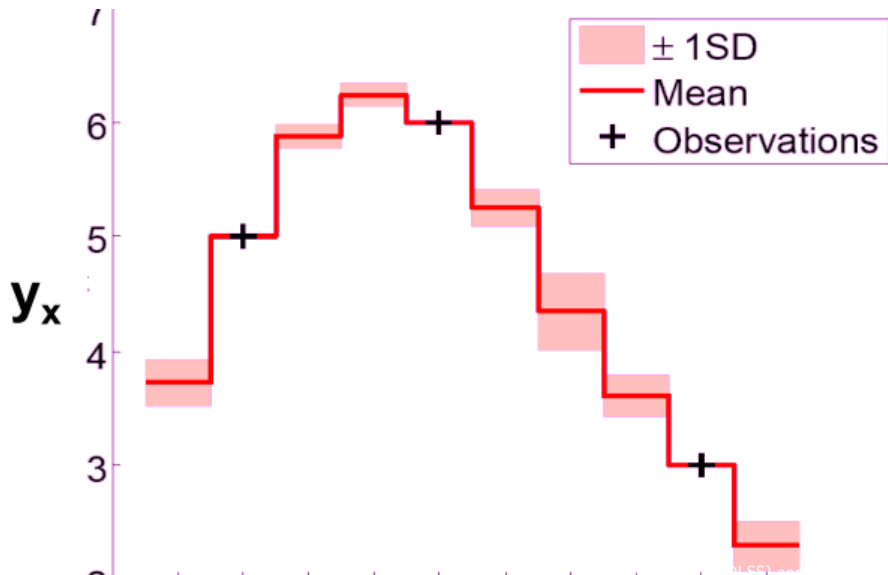
Allows for distributions for variables conditioned on any other observed variables



- A Gaussian process is the generalization of a multivariate Gaussian distribution to a potentially infinite number of variables
- Gives us the limit of potentially infinite number of variables infinitesimally closer together represented by an infinite-length dimension vector
- Provides non-parametric model for functions defined by mean and covariance

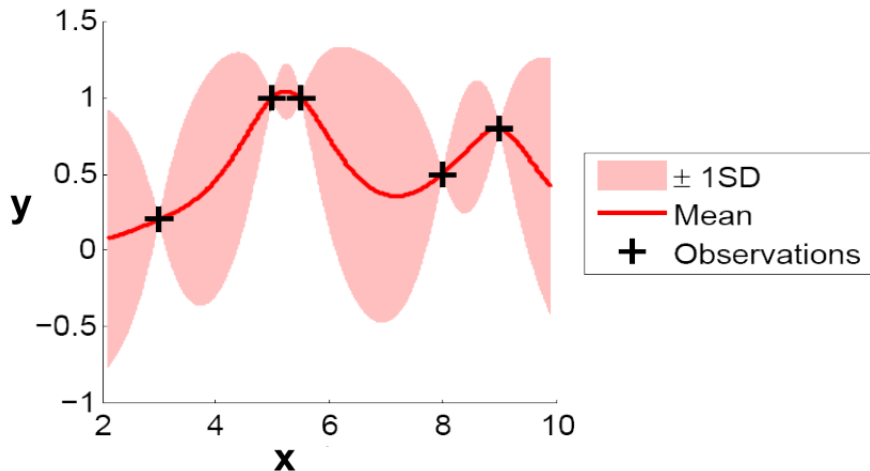
Gaussian Process

Infinite number of variables



Gaussian Process

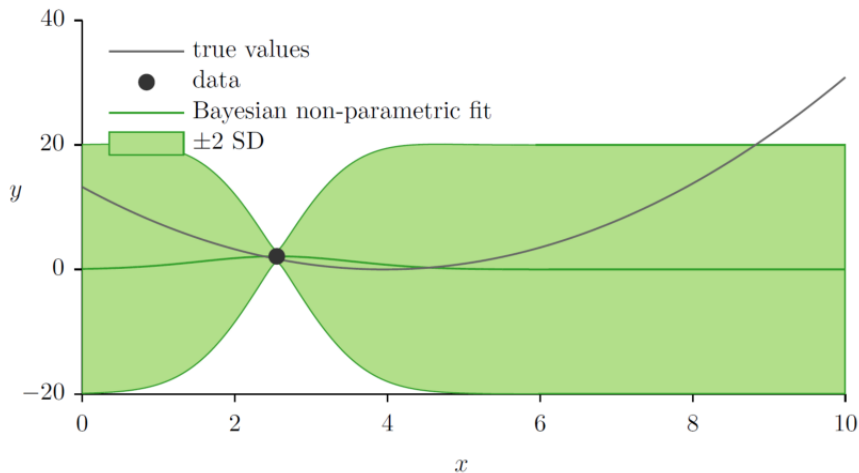
Non-parametric model



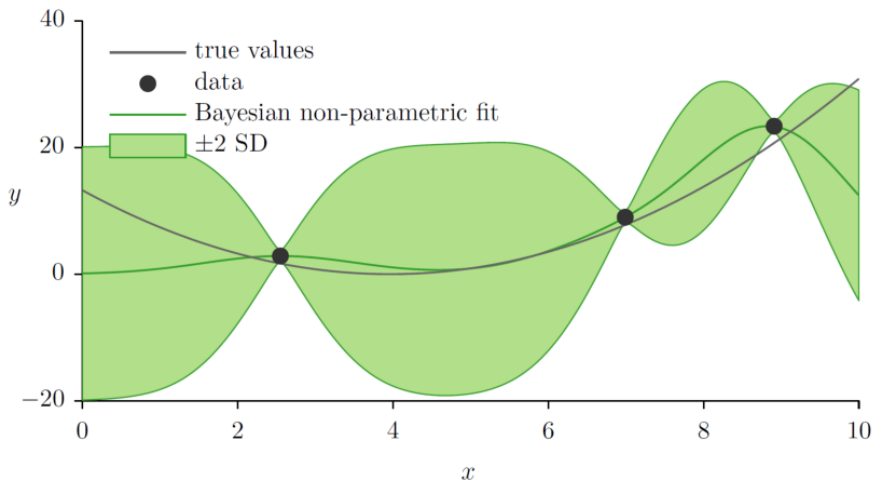
Gaussian Process

- Complexity that grows with data
- Robust to overfitting

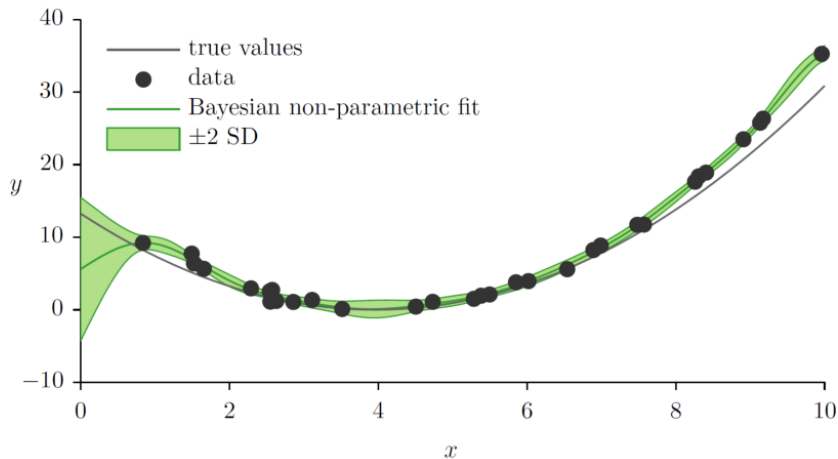
Gaussian Process



Gaussian Process



Gaussian Process



Outline

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- **Bayesian optimization as decision theory**
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- Integration beats Optimization

3 Conclusion

- Experiments

Bayesian Optimization

- Bayesian optimization is the approach of probabilistically modelling $f(x,y)$ and using decision theory to make optimal use of computation
- by defining the costs of observation and uncertainty, we can select evaluations optimally by minimizing the expected loss with respect to a probability distribution
- Representing the core components: cost evaluation and degree of uncertainty

- loss function - lowest function value found after algorithm ends
- Take a myopic approximation and consider only the next evaluation
- The expected loss is the expected lowest value of the function we've evaluated after the next iteration

Myopic Loss

Consider only with one evaluation remaining, the loss of returning value y with current lowest value μ

$$\lambda(y) \triangleq \begin{cases} y; & y < \eta \\ \eta; & y \geq \eta \end{cases}.$$

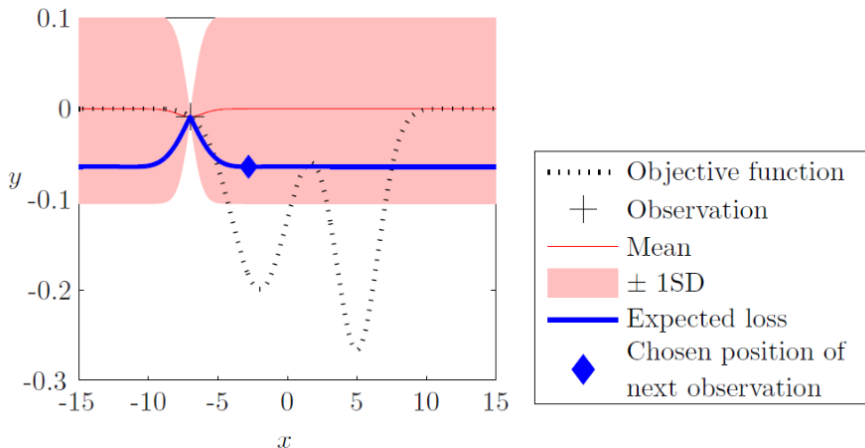
Expected Loss

Expected loss is the expected lowest value

$$\int \lambda(y) p(y \mid x, I_0) \, dy$$

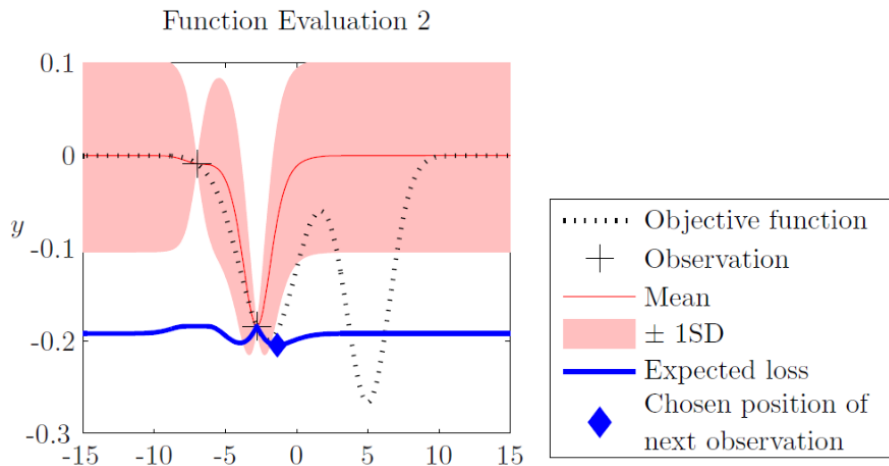
Expected Loss

Use a Gaussian process as the probability distribution for the objective function



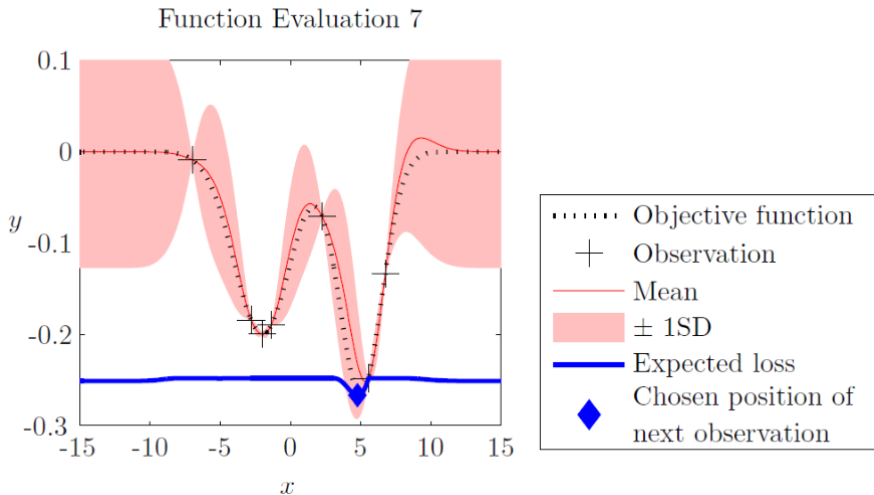
Expected Loss

Exploitative step



Expected Loss

Exploratory step



Outline

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- **Bayesian optimization for tuning hyperparameters**
- Bayesian stochastic optimization
- Integration beats Optimization

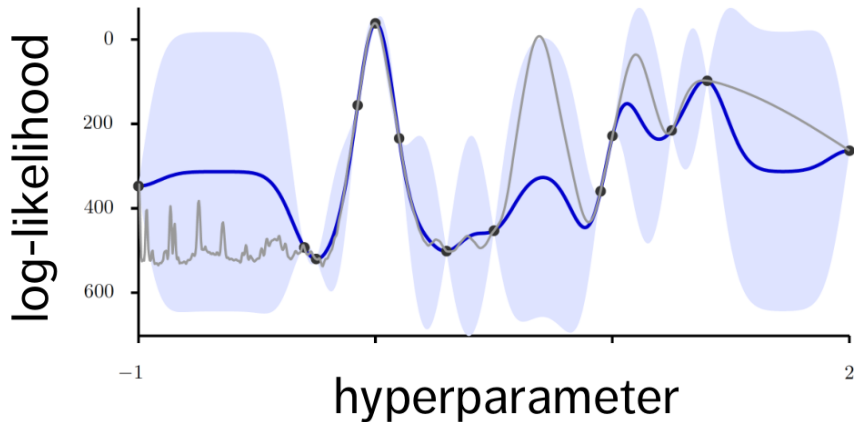
3 Conclusion

- Experiments

- Tuning to cope with model parameters like periods
- Optimization gives a reasonable heuristic
- But Bayesian optimization better

Bayesian Optimization

Better representation across hyperparameters

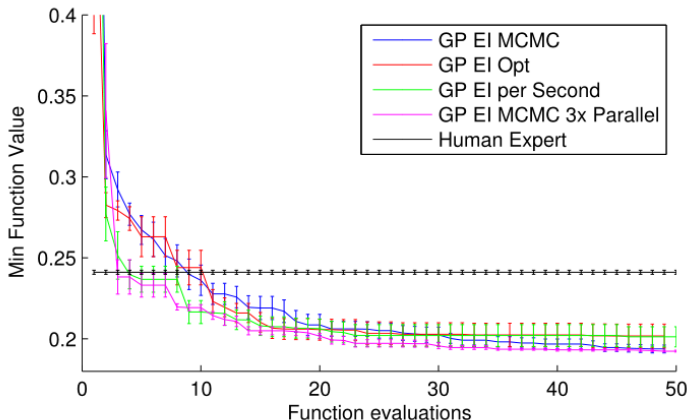


Bayesian Optimization

Tune convolutional neural networks

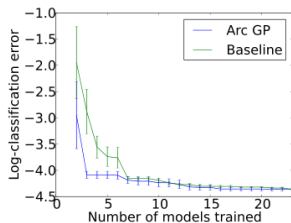
Allows defining the right prior information

Snoek, Larochelle and Adams (2012)

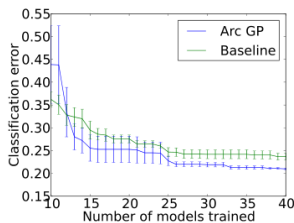


Bayesian Optimization

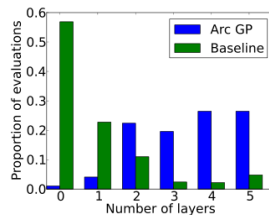
Automated structure learning Swersky et al (2013)



(a) MNIST



(b) CIFAR-10



(c) Architectures searched

Outline

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- **Bayesian stochastic optimization**
- Integration beats Optimization

3 Conclusion

- Experiments

- Using only a subset of the data gives a noisy likelihood evaluation
- Use Bayesian optimization for stochastic learning

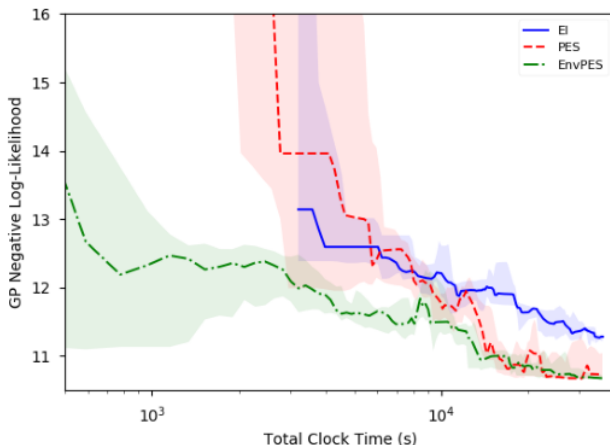
- Within Bayesian Optimization noise is not a problem
- If additional noise in the random variable we can just add a noise likelihood to complement model
- Encode that cost as a function of the number of data
- Intelligently choose the size of data that it needs at runtime to best optimization

Bayesian Optimization

Batch size

Klein, Falkner, Bartels, Hennig, Hutter (2017);

McLeod, Osborne Roberts (2017)



A Random Number

- is epistemic (personal particular to an agent) (computation is always conditional on prior knowledge)
- use useful to foil a malicious adversary (few in numerics)
- is never the minimizer of an expected loss (only when totally flat)

Outline

1 Introduction

- Probabilistic Numerics

2 Components

- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- **Integration beats Optimization**

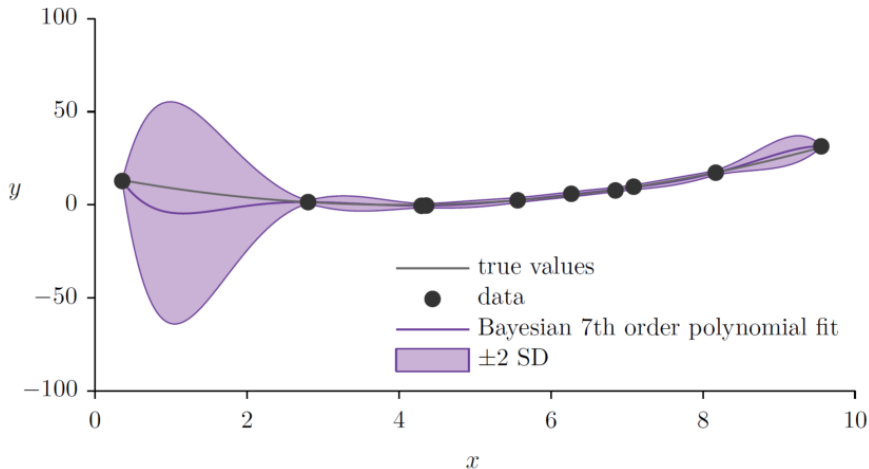
3 Conclusion

- Experiments

- Naive fitting can lead to overfitting

Integrating

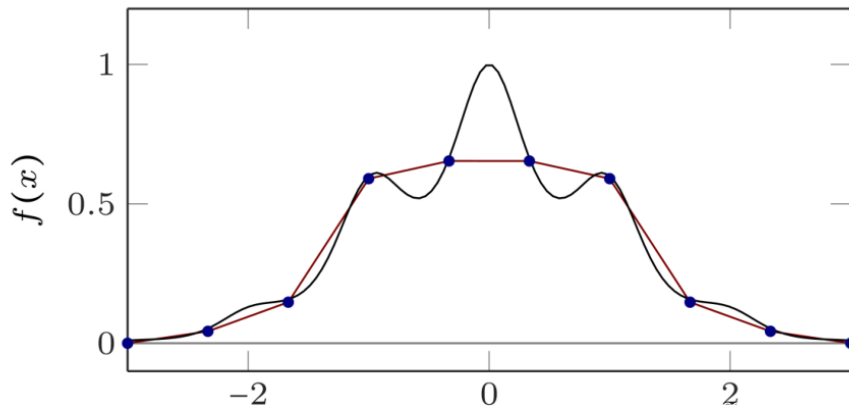
Reduces overfitting and estimates uncertainty



Integrating

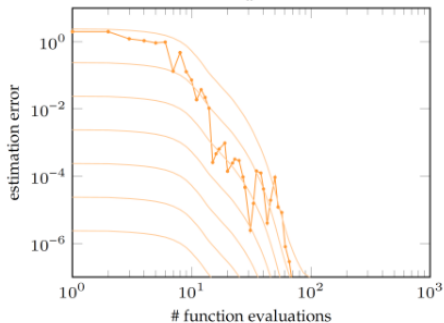
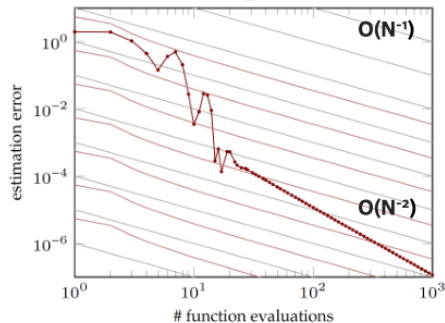
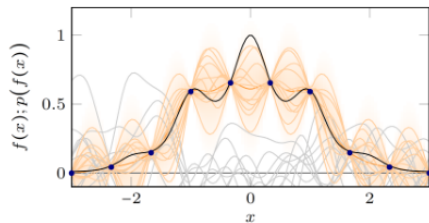
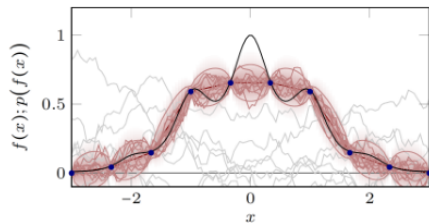
Don't average use quadrature

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



Bayesian Quadrature

Trapezoid method



Outline

1 Introduction

- Probabilistic Numerics

2 Components

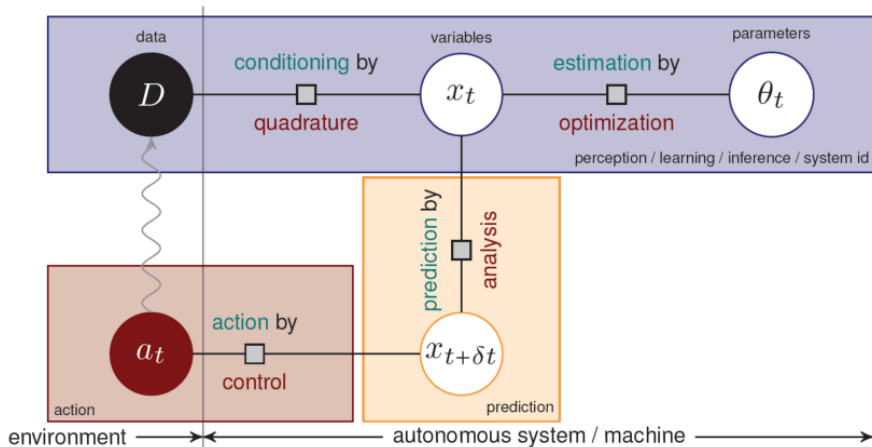
- Probabilistic modeling of functions
- Bayesian optimization as decision theory
- Bayesian optimization for tuning hyperparameters
- Bayesian stochastic optimization
- Integration beats Optimization

3 Conclusion

- Experiments

Model

Propagates uncertainty



Model

Converges

synthetic (moG)

