

# Gaussian random fields in machine learning

Viacheslav Borovitskiy (Slava)

St. Petersburg State University  
St. Petersburg Department of Steklov Mathematical Institute

19 November 2020

# Talk structure

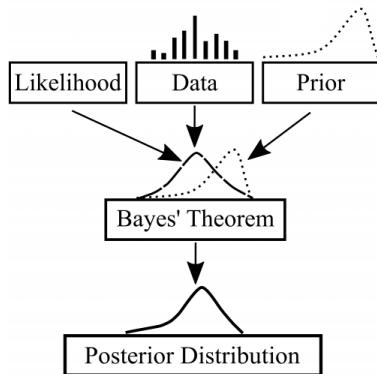
- 1 Introduction
- 2 Gaussian processes
- 3 Applications
- 4 Our own research

# Talk structure

- 1 Introduction
- 2 Gaussian processes
- 3 Applications
- 4 Our own research

# Gaussian processes in machine learning

Bayesian learning paradigm:



Gaussian processes (GPs) — non-parametric prior over functions.

GPs are indeed useful

---

# Bayesian Optimization in AlphaGo

---

**Yutian Chen, Aja Huang, Ziyu Wang, Ioannis Antonoglou, Julian Schrittwieser,  
David Silver & Nando de Freitas**

DeepMind, London, UK  
yutianc@google.com

They used GPs to model target function and guide decision (optimization) process.

# Talk structure

- 1 Introduction
- 2 Gaussian processes
- 3 Applications
- 4 Our own research

# Gaussian process regression

GP — distribution over functions.

Bayesian inference for GPs:

- prior:** hand-picked GP

- data:** noisy evaluations of the function

- likelihood:** induced by Gaussian noise assumption

- posterior:** another GP

Let us explore this visually ...

# Visual guide to Gaussian process regression





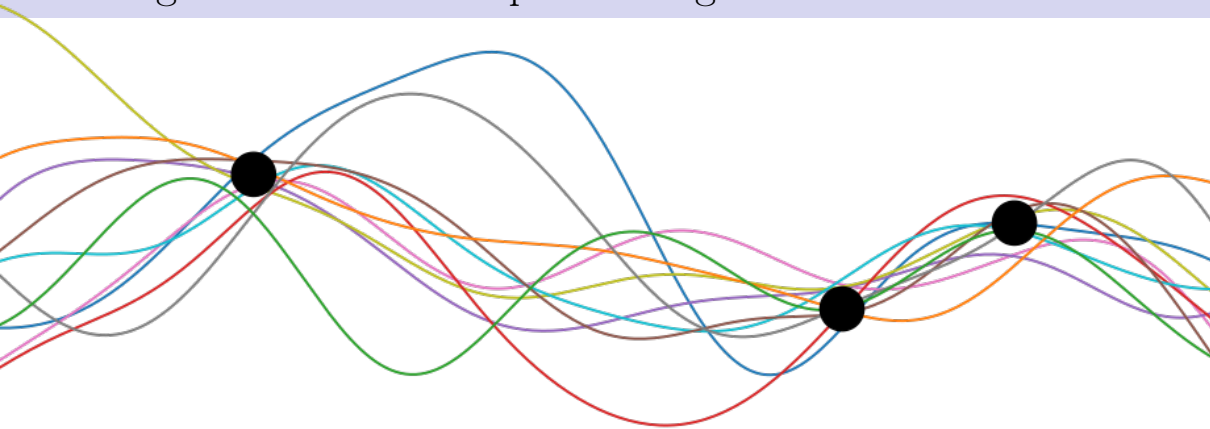
# Visual guide to Gaussian process regression



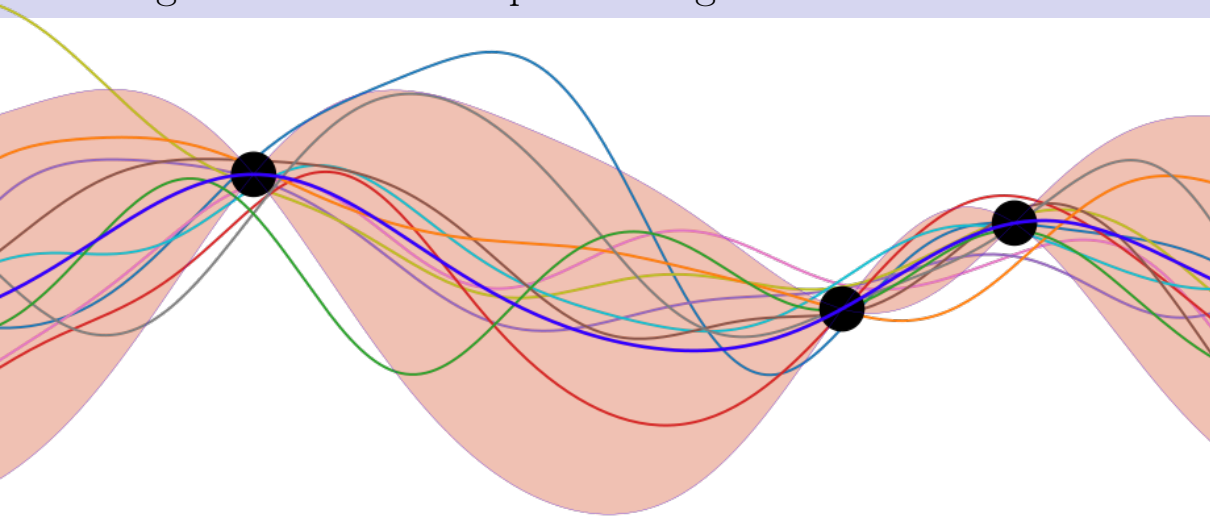
# Visual guide to Gaussian process regression



# Visual guide to Gaussian process regression



# Visual guide to Gaussian process regression



# What is a Gaussian process?

## Gaussian random variable

- distribution over  $\mathbb{R}$ , denoted by  $N(\mu, \sigma^2)$ ,
- determined by two numbers: mean  $\mu$  and variance  $\sigma^2$ .

## Multivariate Gaussian random variable

- distribution over  $\mathbb{R}^d$ , denoted by  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,
- determined by the mean vector  $\boldsymbol{\mu}$  and the covariance matrix  $\boldsymbol{\Sigma}$ .

## Gaussian process

- distribution over functions from  $X$  to  $\mathbb{R}$ , denoted by  $GP(m, k)$ ,
- determined by two functions  $m : X \rightarrow \mathbb{R}$  (mean) and  $k : X \times X \rightarrow \mathbb{R}$  (covariance).

Gaussian processes are appealing in practice due to their simplicity (among other stochastic processes).

# Bayesian inference for GPs

Bayesian inference for GPs takes in

- a prior distribution over functions of form  $GP(m, k)$ ,
- noisy evaluations  $y_1, \dots, y_n$  of the unknown function of interest at  $x_1, \dots, x_n$ .

and returns the distribution over functions of form

$$GP(\tilde{m}, \tilde{k}).$$

Given  $m$  and  $k$ , the functions  $\tilde{m}$  and  $\tilde{k}$  can be computed in a finite time. Specifically:

$$\begin{aligned}\tilde{m}(u) &= m(u) + \mathbf{K}_{f(u)f(x)} \left( \mathbf{K}_{f(x)f(x)} + \sigma^2 I \right)^{-1} (\mathbf{y} - m(\mathbf{x})) \\ \tilde{k}(u, v) &= k(u, v) - \underbrace{\mathbf{K}_{f(u)f(x)}}_{\text{vector } 1 \times n} \underbrace{\left( \mathbf{K}_{f(x)f(x)} + \sigma^2 I \right)^{-1}}_{\text{matrix } n \times n} \underbrace{\mathbf{K}_{f(x)f(v)}}_{\text{vector } n \times 1}.\end{aligned}$$

# The Gaussian process regression algorithm

So how do we turn the data  $(x_1, y_1), \dots, (x_n, y_n)$  into a reasonable stochastic model interpolating it?

- ❶ Come up with a parametric families  $m_\theta$  and  $k_\theta$  for prior mean and covariance functions.
- ❷ Use maximum likelihood estimation to pick the optimal set of parameters  $\theta$  and the optimal noise value  $\sigma^2$  from data  $(x_1, y_1), \dots, (x_n, y_n)$ .
- ❸ Perform Bayesian inference with prior  $GP(m_\theta, k_\theta)$ , data  $(x_1, y_1), \dots, (x_n, y_n)$  and likelihood noise  $\sigma^2$ .

As a result, obtain the posterior  $\tilde{m}$  and  $\tilde{k}$ .

- ❹ Use
  - ▶  $N(\tilde{m}(u), \tilde{k}(u, u))$  as a stochastic prognosis at a new location  $u$ .
  - ▶ use samples of  $GP(\tilde{m}, \tilde{k})$  as an ensemble of possible deterministic models.



# Visual guide to Gaussian process regression





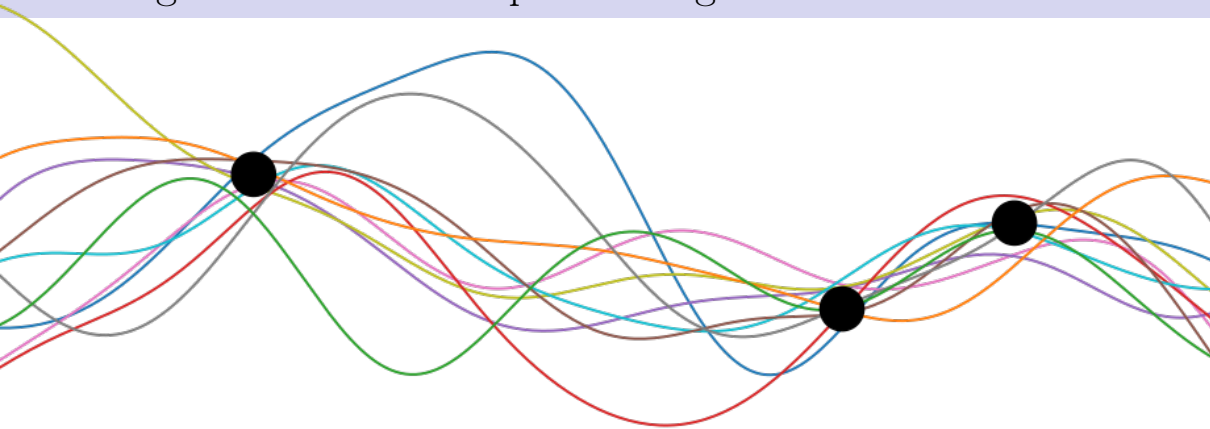
# Visual guide to Gaussian process regression



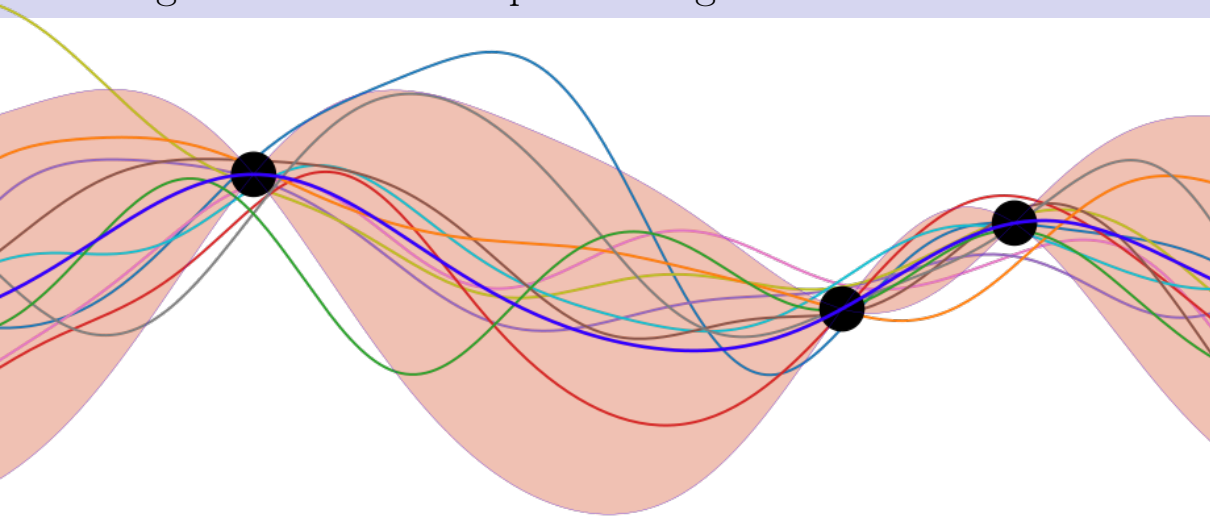
# Visual guide to Gaussian process regression



# Visual guide to Gaussian process regression



# Visual guide to Gaussian process regression



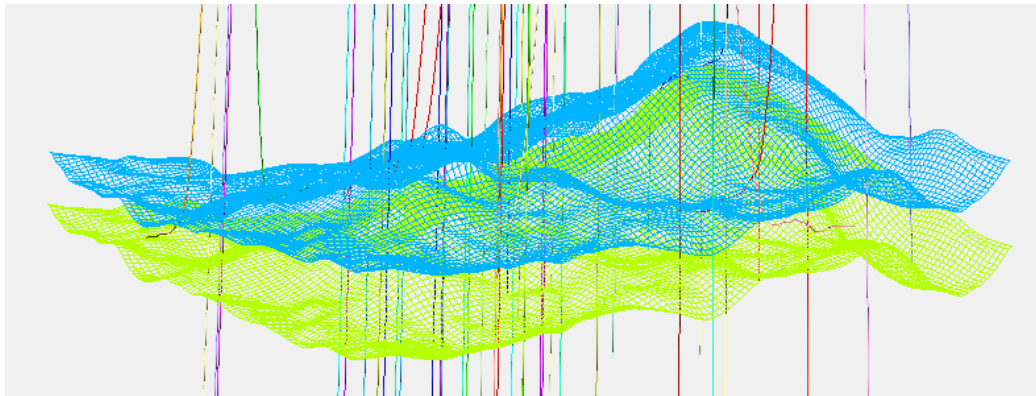
# Talk structure

- 1 Introduction
- 2 Gaussian processes
- 3 Applications**
- 4 Our own research

# Geostatistical modeling of petroleum reservoirs

Problem: interpolate well data into the interwell space.

The data is very sparse, thus deterministic model is undesirable.



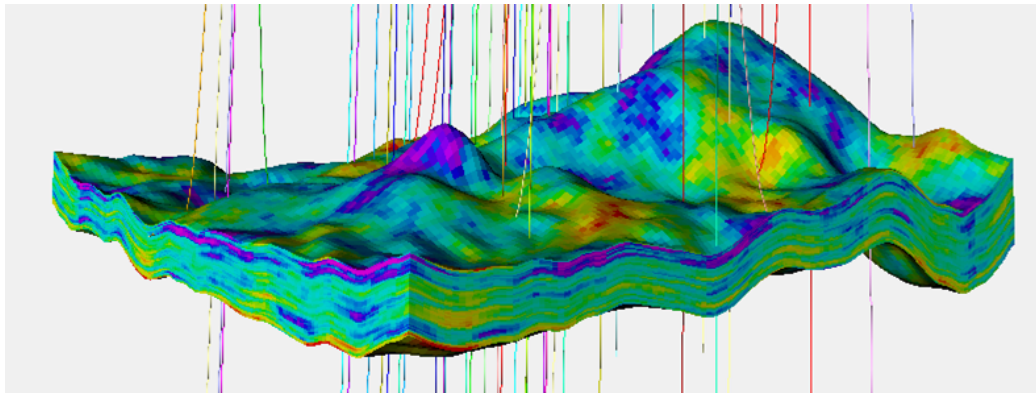
Reservoir structure, well locations.



# Geostatistical modeling of petroleum reservoirs

Problem: interpolate well data into the interwell space.

The data is very sparse, thus deterministic model is undesirable.



A single sample of a Gaussian process model in the interwell space

# Bayesian optimization of expensive black-box functions

Problem: minimize the target function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ .

At  $n$ 'th step  $\phi$  has already been evaluated at  $x_1, \dots, x_n$ . How do we choose  $x_{n+1}$ ?

Build posterior GP  $f$  using data

$$x_1, \dots, x_n, \quad \phi(x_1), \dots, \phi(x_n).$$

Choose

$$x_{n+1} = \arg \max_{x \in \mathbb{R}^d} \mathbb{P}(f(x) < \min_{i=1..n} \phi(x_i)). \quad (MPI)$$

or

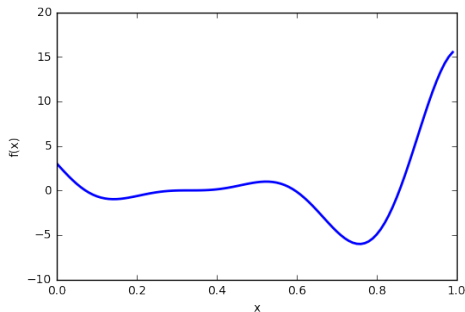
$$x_{n+1} = \arg \max_{x \in \mathbb{R}^d} \mathbb{E} \max(\min_{i=1..n} \phi(x_n) - f(x), 0). \quad (EI)$$

Automatic exploration/exploitation trade-off.



# Example

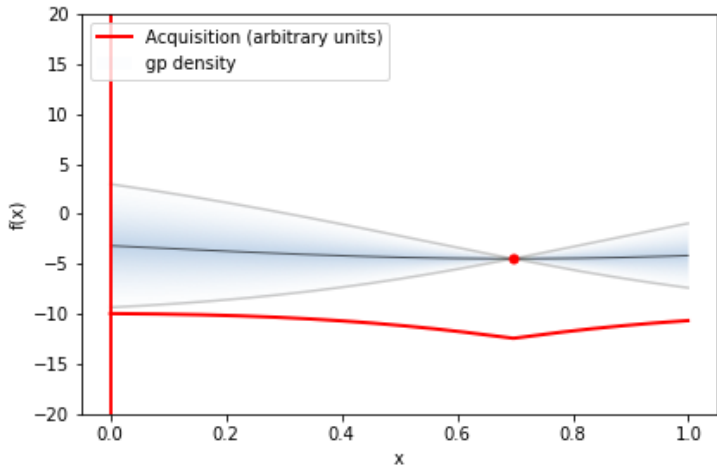
Let us minimize Forrester function  $f(x) = (6x - 2)^2 \sin(12x - 4)$ .



Choose some prior as  $f_0 \sim \text{GP}(?, ?)$ .

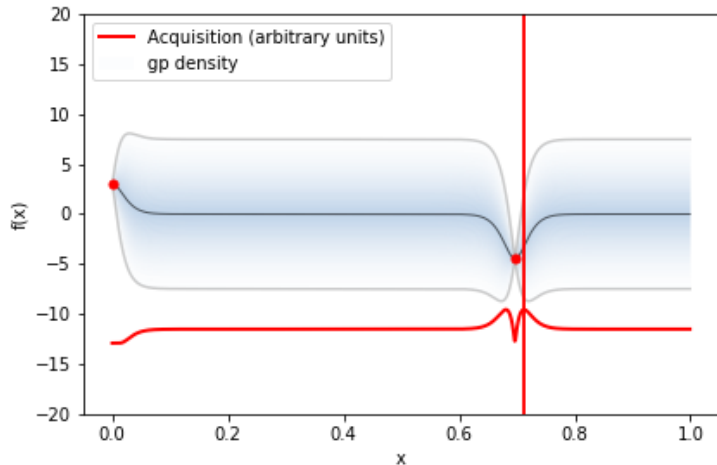
# Example

Iteration 1.



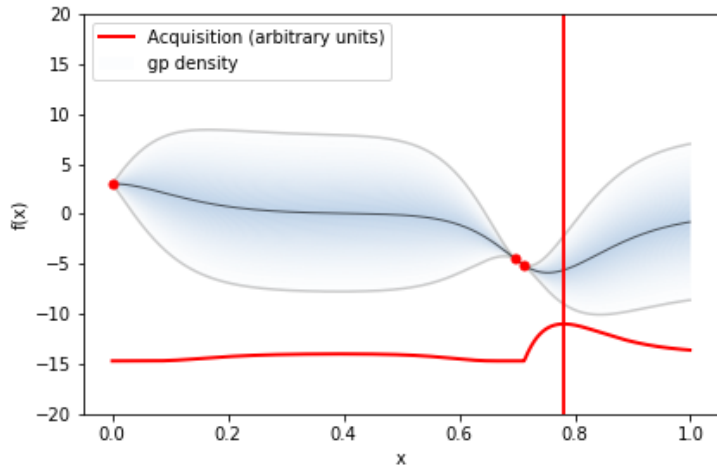
# Example

Iteration 2.



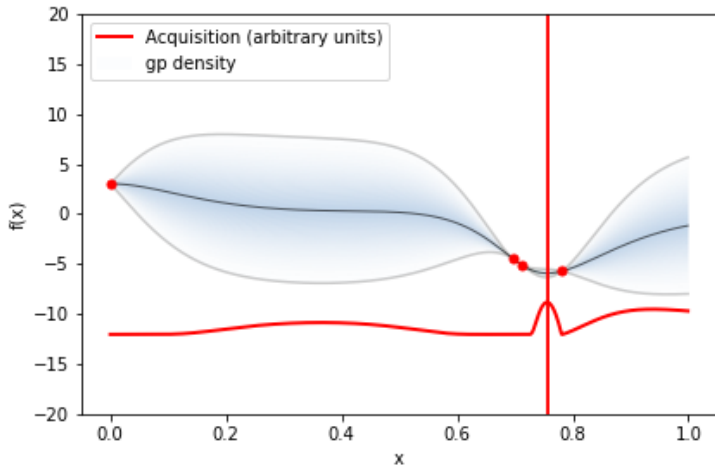
# Example

Iteration 3.



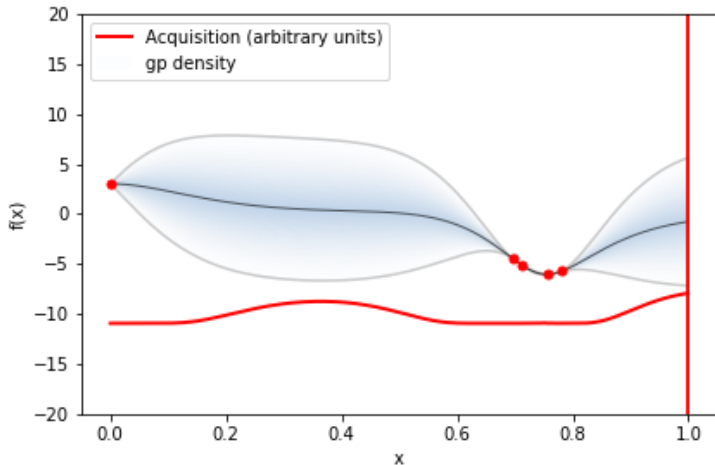
# Example

Iteration 4.



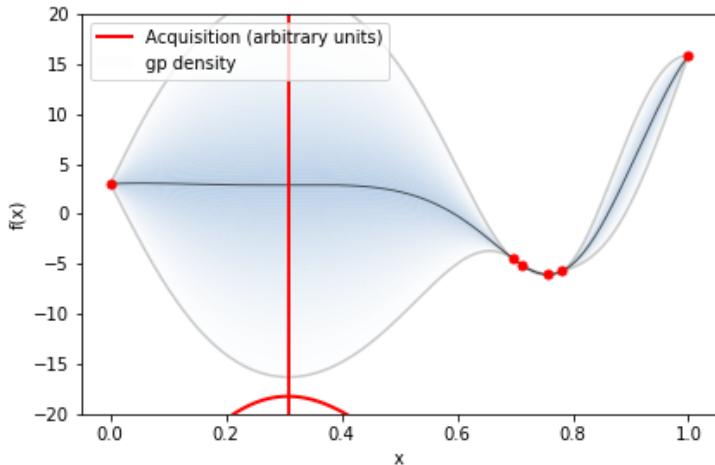
# Example

Iteration 5.



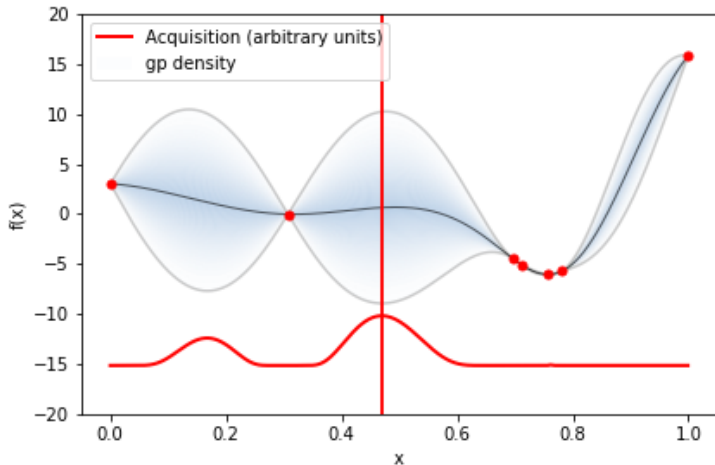
# Example

Iteration 6.



# Example

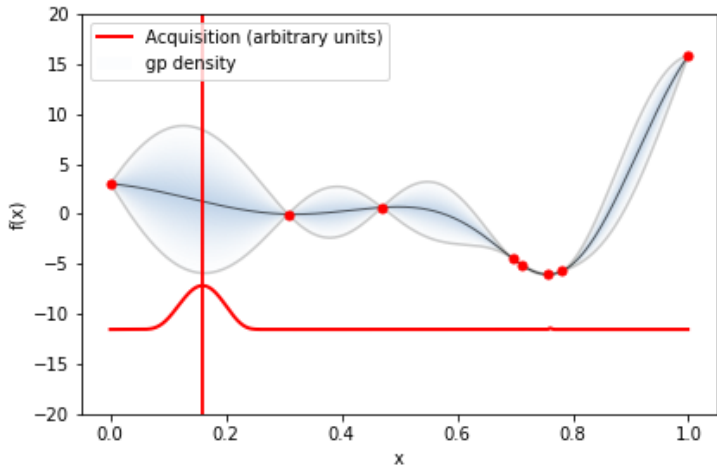
Iteration 7.





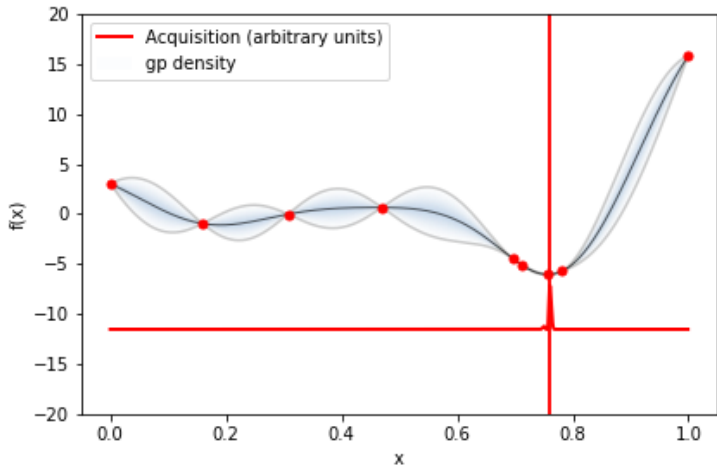
# Example

Iteration 8.



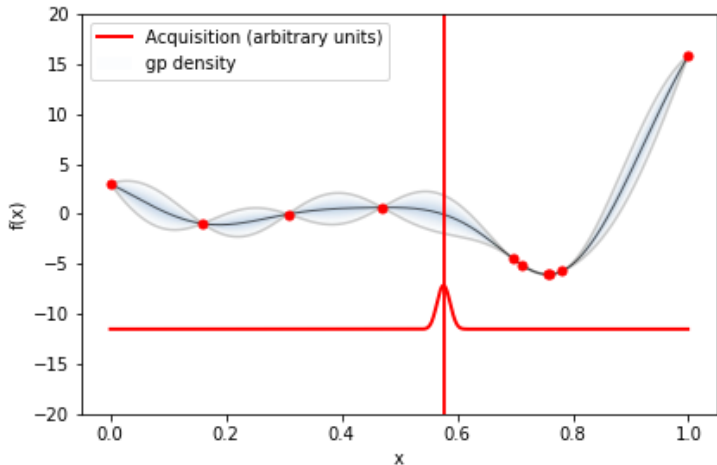
# Example

Iteration 9.



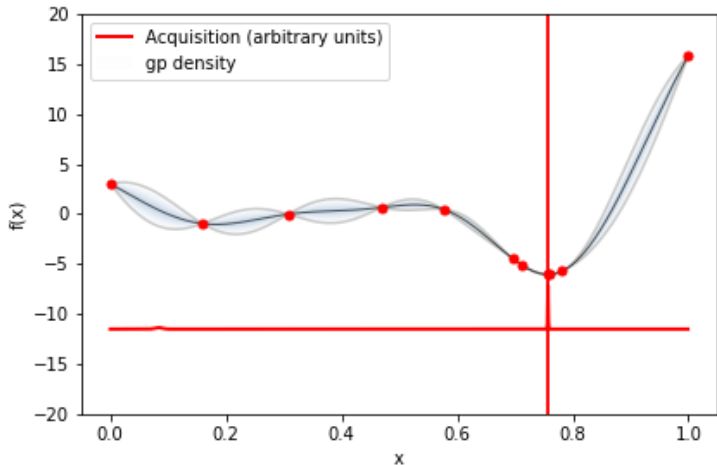
# Example

Iteration 10.



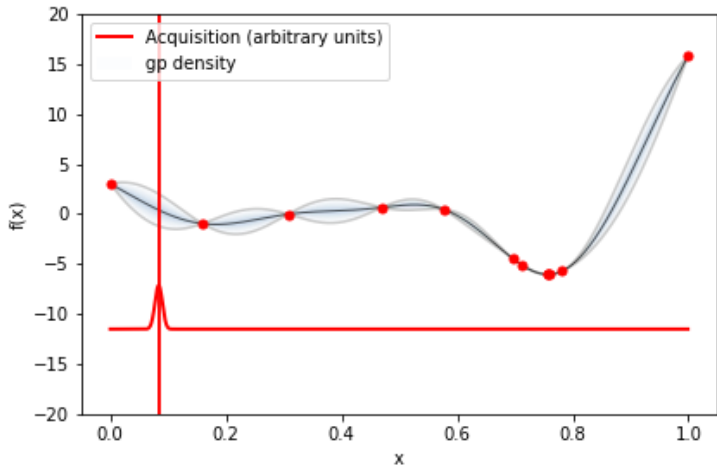
# Example

Iteration 11.



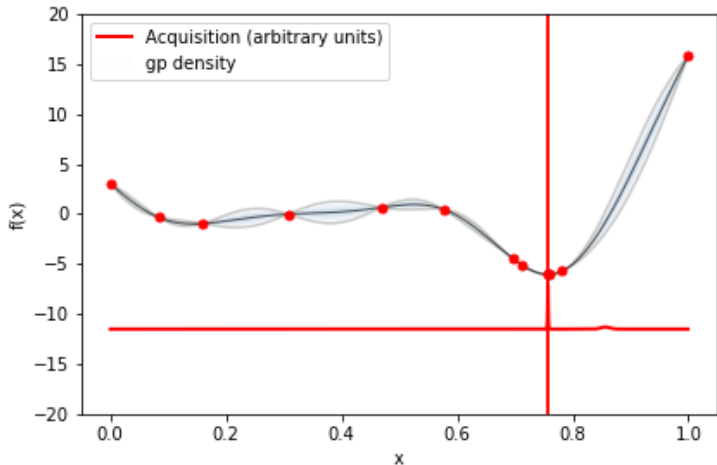
# Example

Iteration 12.



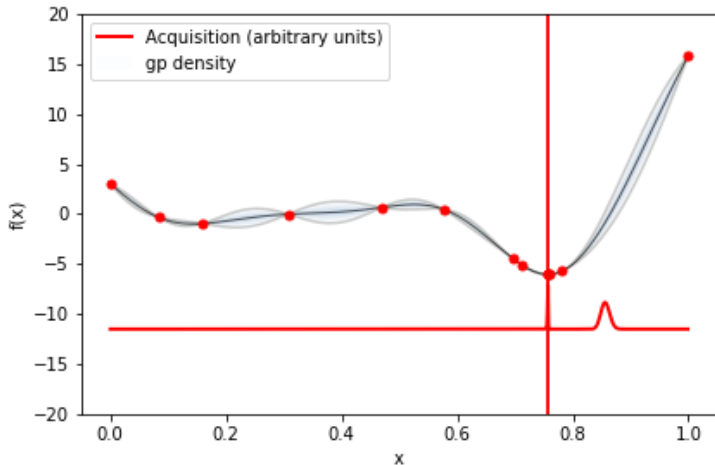
# Example

Iteration 13.



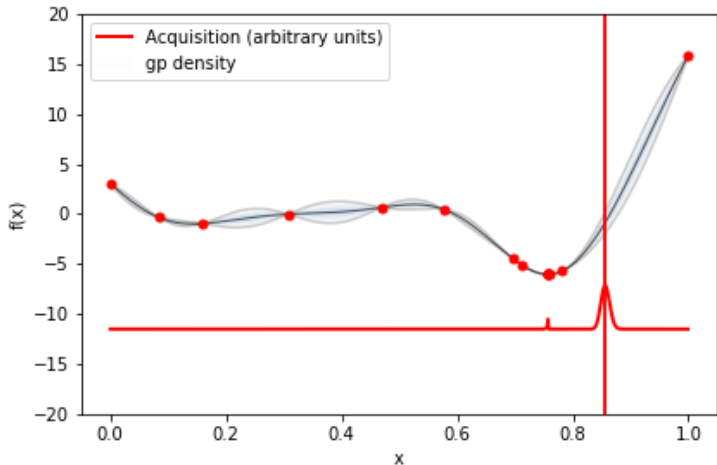
# Example

Iteration 14.



# Example

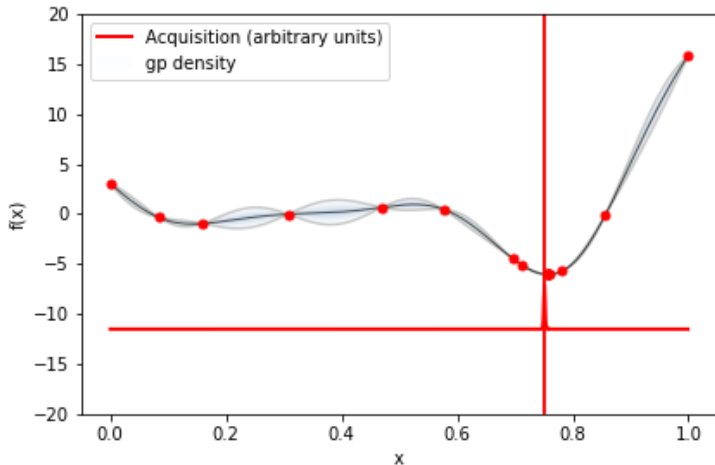
Iteration 15.





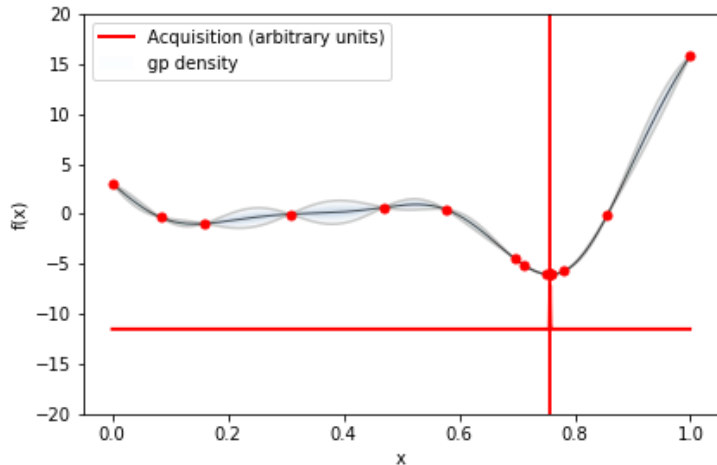
# Example

Iteration 16.



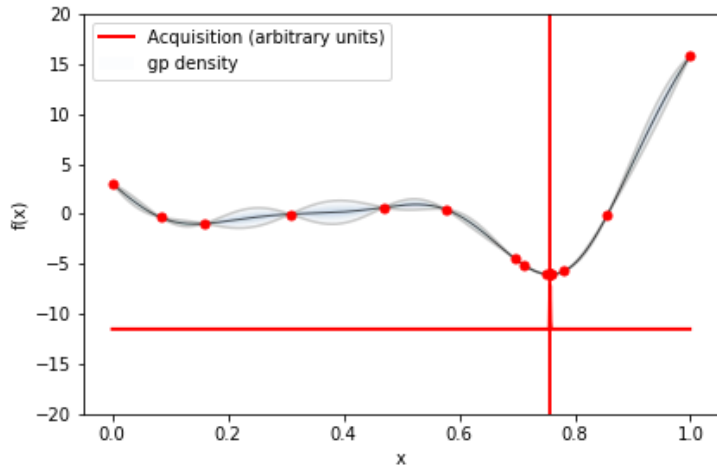
# Example

Iteration 17.



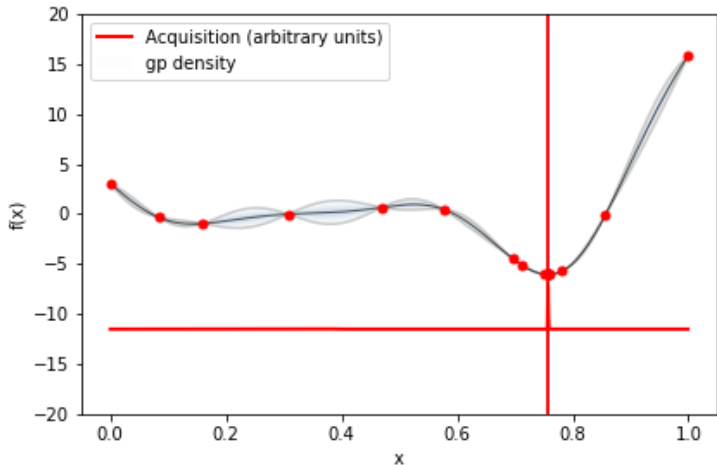
# Example

Iteration 18.



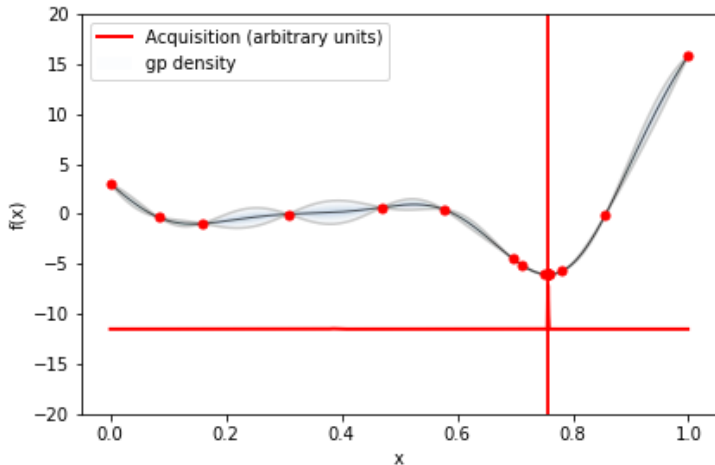
# Example

Iteration 19.



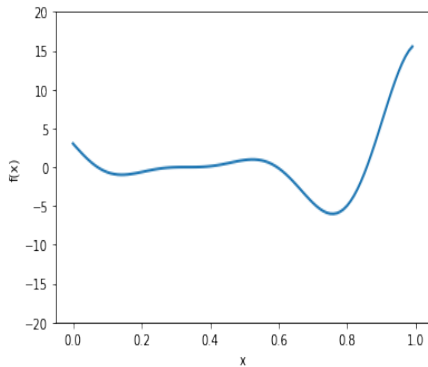
# Example

Iteration 20.

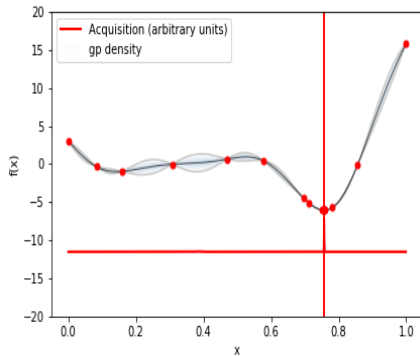


# Example

Let us compare the model after 20 iterations with the target function.



(a) Target function



(b) GP model

Classical control problem: physics is known, find optimal control.

Reinforcement learning control problem: physics is unknown, try to learn physics from data and on the go build the optimal control.

Second approach is supposed to bring us the cheap robots, for which

- we don't indeed know the physics (it deviates too much from the “ideal”),
- learning this physics by hand is of course possible, but it increases the price.

# PILCO for robotics and control

PILCO (Probabilistic Inference for Learning COntrol) — an approach that uses GPs to model the unknown physics.

---

## PILCO: A Model-Based and Data-Efficient Approach to Policy Search

---

**Marc Peter Deisenroth**

Department of Computer Science & Engineering, University of Washington, USA

MARC@CS.WASHINGTON.EDU

**Carl Edward Rasmussen**

Department of Engineering, University of Cambridge, UK

CER54@CAM.AC.UK

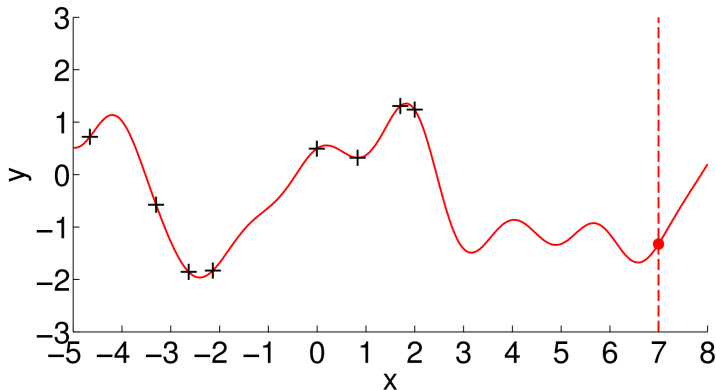
The model can be described by  $x_{t+1} = f(x_t, u_t) + w$ , where

- $x_t$  — trajectory,
- $u_t$  — control,
- $f$  models physics,
- $w \sim N(0, \sigma^2)$  — random noise.



# PILCO for robotics and control

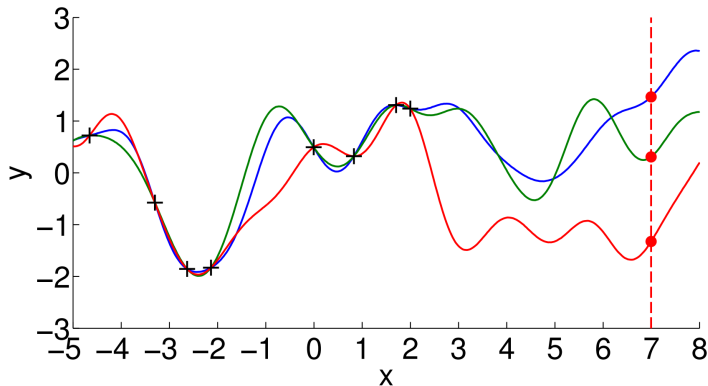
Imagine that  $f$  is modeled deterministically.



Consider a prognosis at  $x = 7$ .

# PILCO for robotics and control

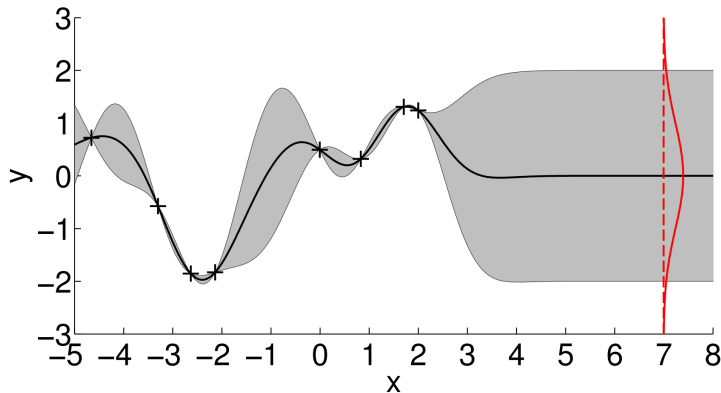
Imagine that  $f$  is modeled deterministically.



There exists a number of plausible models and thus a number of different predictions.

# PILCO for robotics and control

What if we model  $f$  as a GP?

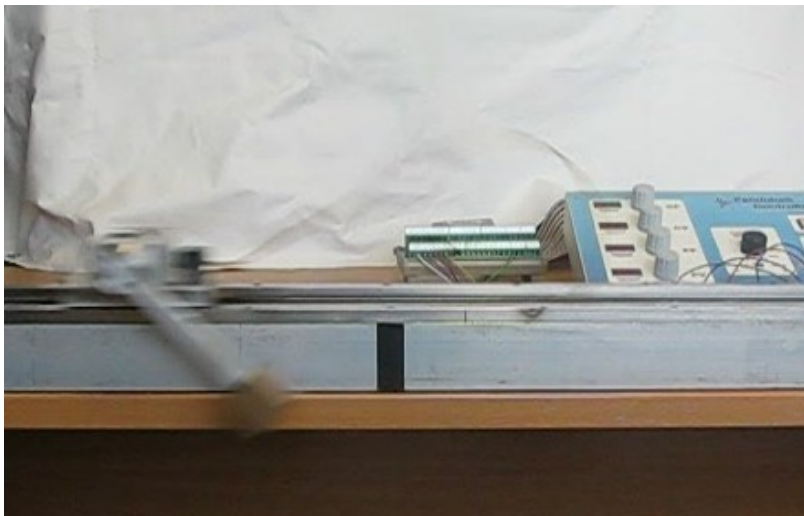


If we use GPs, we are able to use an infinite number of plausible models all at once.

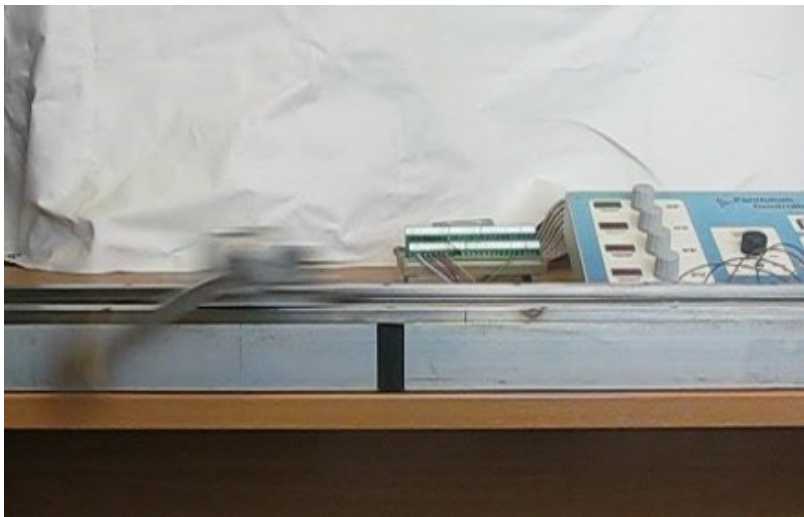
# Example: learning to control a pendulum



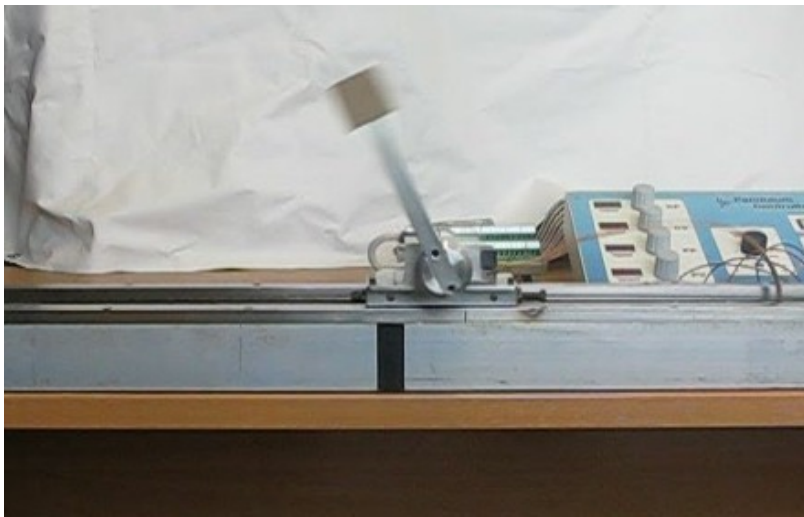
# Example: learning to control a pendulum



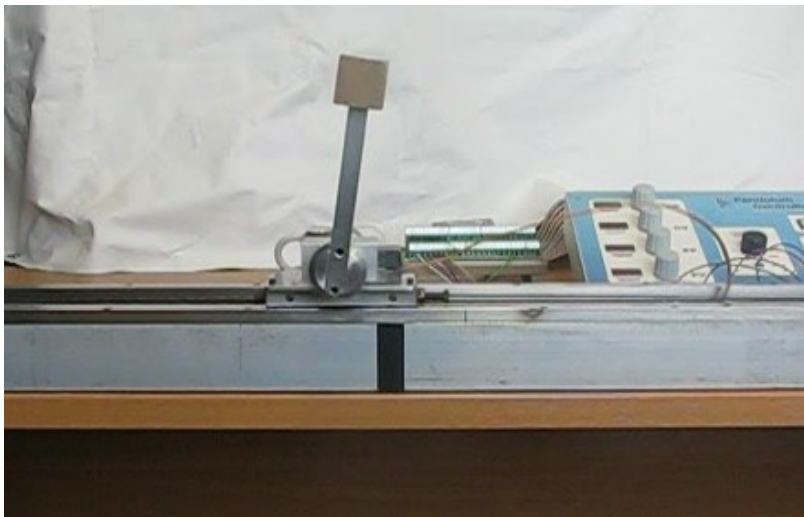
# Example: learning to control a pendulum



# Example: learning to control a pendulum

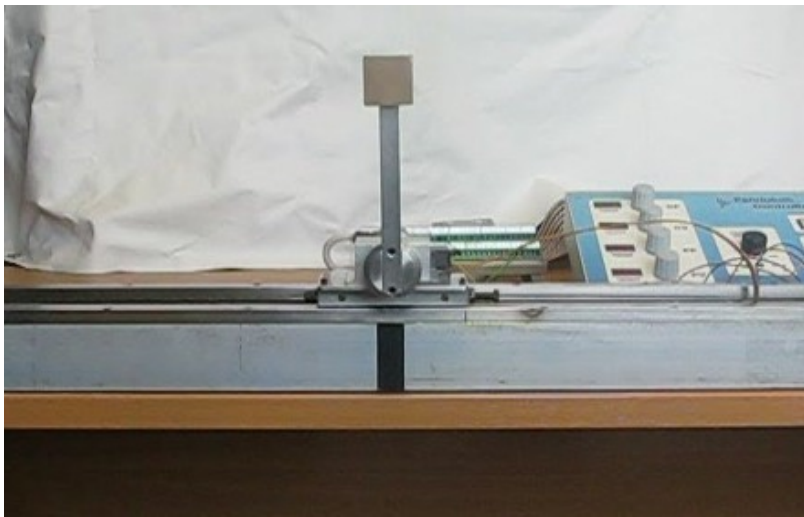


# Example: learning to control a pendulum





# Example: learning to control a pendulum



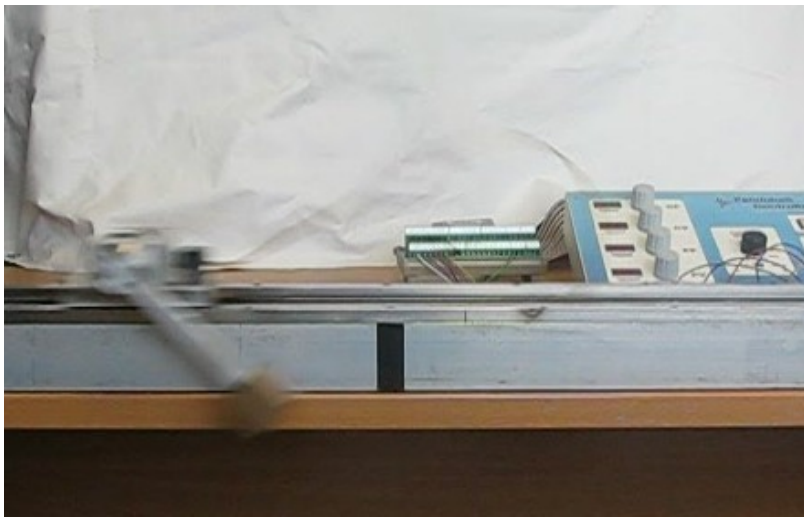
# Example: learning to control a pendulum

Once more...

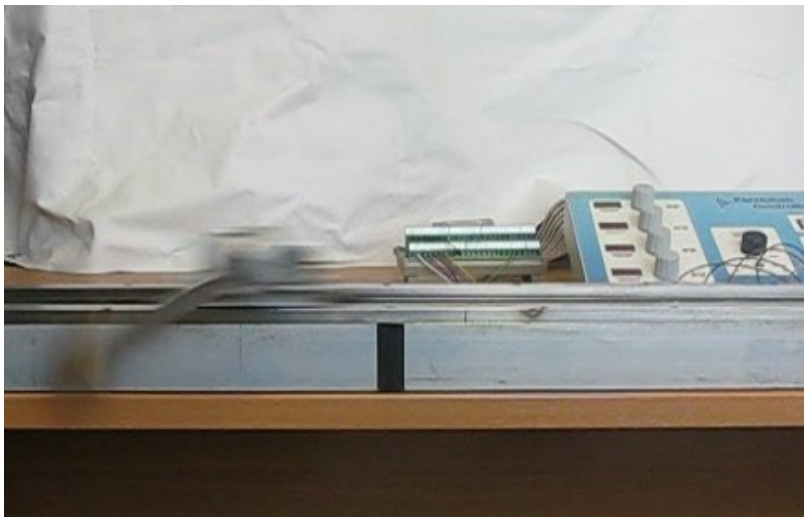
# Example: learning to control a pendulum



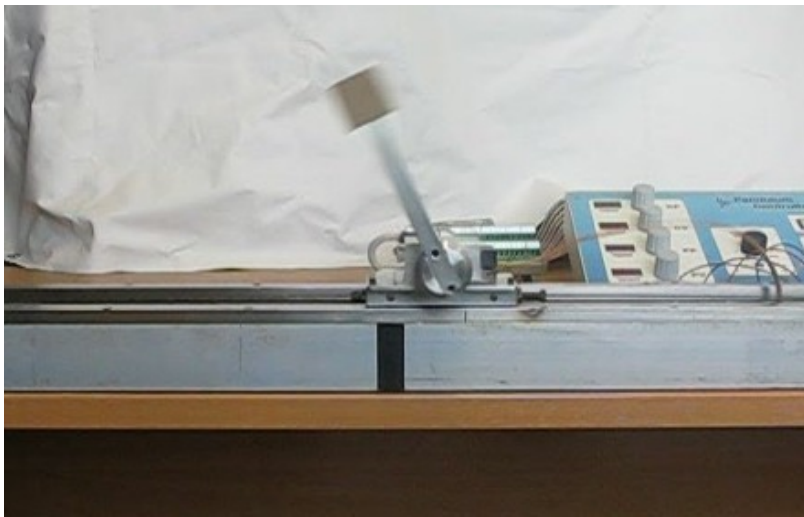
# Example: learning to control a pendulum



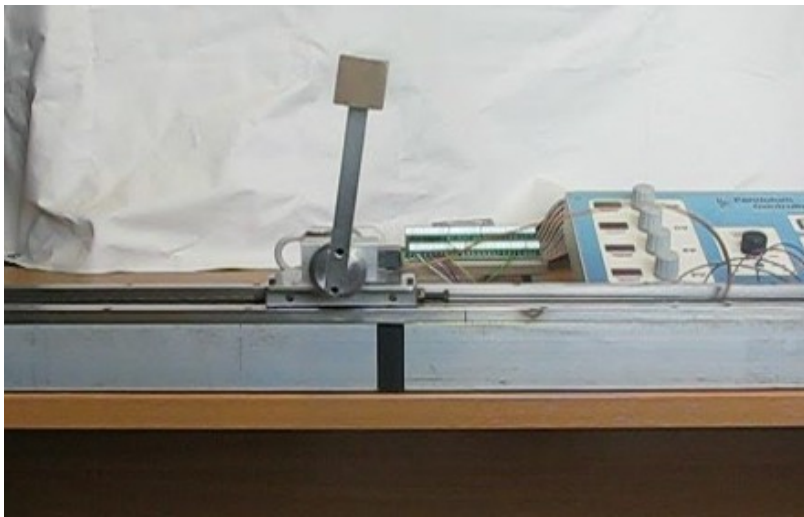
# Example: learning to control a pendulum



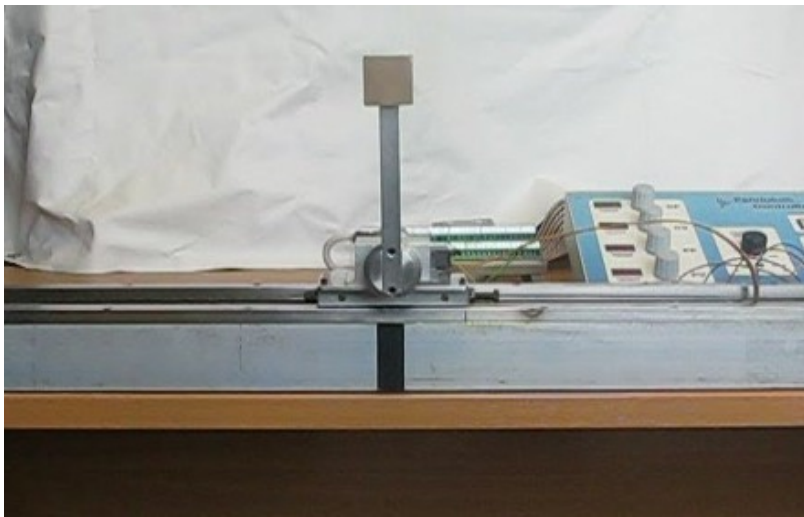
# Example: learning to control a pendulum



# Example: learning to control a pendulum



# Example: learning to control a pendulum





# Talk structure

- 1 Introduction
- 2 Gaussian processes
- 3 Applications
- 4 Our own research

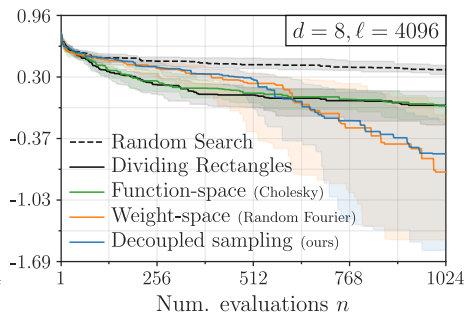
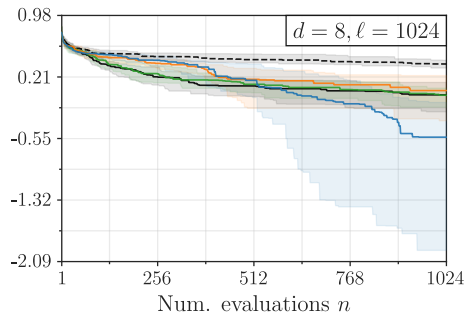
# Efficient sampling from (approximate) posteriors

## Efficiently Sampling Functions from Gaussian Process Posteriors

James T. Wilson<sup>\*1</sup> Viacheslav Borovitskiy<sup>\*2,3</sup> Alexander Terenin<sup>\*1</sup>  
Peter Mostowsky<sup>\*2</sup> Marc Peter Deisenroth<sup>4</sup>



Outstanding  
Paper Honorable  
Mention Award  
at ICML 2020



✓ Improved performance owing to smaller error

## Matérn Gaussian processes on Riemannian manifolds

Viacheslav Borovitskiy<sup>\*1,4</sup> Alexander Terenin<sup>\*2</sup> Peter Mostowsky<sup>\*1</sup> Marc Peter Deisenroth<sup>3</sup>

To be presented on NeurIPS 2020.



(a) Ground truth



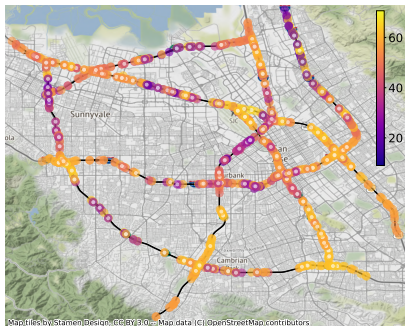
(b) Posterior mean



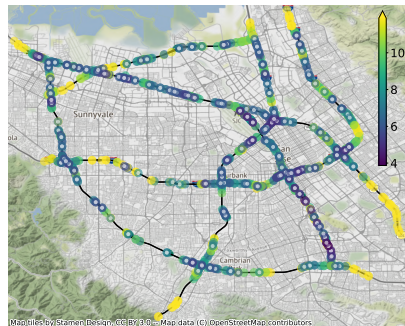
(c) Standard deviation

## Matérn Gaussian Processes on Graphs

In review for AISTATS 2021



(a) Mean



(b) Standard deviation

# Thank you for your attention!

viacheslav.borovitskiy@gmail.com



St Petersburg  
University

Mathematics & Computer Science department

Some figures were taken from: <http://inverseprobability.com/talks/>.