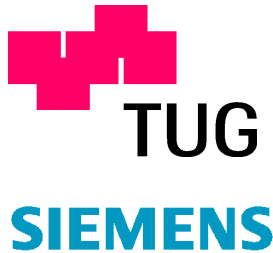


# Some Aspects of Learning with Gaussian Processes

Anton Schwaighofer, May 2003



TU Graz, Institute for Theoretical  
Computer Science

Siemens Corporate Technology  
Department of Neural Computation

# Overview

---

- Gaussian processes basics
- Gaussian processes for Bayesian regression (GPR)
- The evidence procedure in GPR
- Choice of kernel functions
- Application areas

# Gaussian Processes: Motivation

We wish to do (Bayesian) inference for predicting the value of some unknown function  $f$  on a test point  $\mathbf{x}^*$ , given noisy measurements (training data)  $\{\mathbf{x}_i, y_i\}_{i=1}^N$

**Neural network** Assume a neural network model with some parameters (weights)  $\mathbf{w}$ , set weights (train the network) to minimize the mean squared error on the training set

**Gaussian processes** Assume a prior distribution over the space of functions that possibly may have generated the training data. Compute the posterior distribution of functions  $p(f | \{\mathbf{x}_i, y_i\})$  given the training data. Use the most likely function to make predictions

Gaussian processes are probably the simplest way of specifying a (non-trivial) prior over function space

# Basics (1)

Stochastic process is a collection of random variables  $\{F(\mathbf{x}) | \mathbf{x} \in X\}$  indexed by a set  $X$

- Signal processing:  $X$  is a time variable
- Here:  $X = \mathbb{R}^D$  (locations in a  $D$ -dimensional space)

Stochastic processes are characterized by the joint distributions of finite subsets  $\{F(\mathbf{x}_1), \dots, F(\mathbf{x}_N)\}$

**Gaussian process (GP):** Any finite subset of  $F$ -variables has a joint multivariate Gaussian distribution. GP is specified by its mean function  $\mu(\mathbf{x})$  (which we assume to be zero) and its covariance function  $k(\mathbf{x}, \mathbf{x}')$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[F(\mathbf{x})F(\mathbf{x}')]$$

## Basics (2)

For any set of locations  $\mathbf{x}_1, \dots, \mathbf{x}_N$  the associated random variables (functional values)  $F(\mathbf{x}_1), \dots, F(\mathbf{x}_N)$  are Gaussian distributed with mean  $\boldsymbol{\mu} = 0$  and covariance matrix  $K$

$$\begin{aligned}(F(\mathbf{x}_1), \dots, F(\mathbf{x}_N)) &\sim N(\mathbf{0}, K) \\ K_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

$k(\cdot, \cdot)$  is the **covariance function (kernel function)** of the GP.

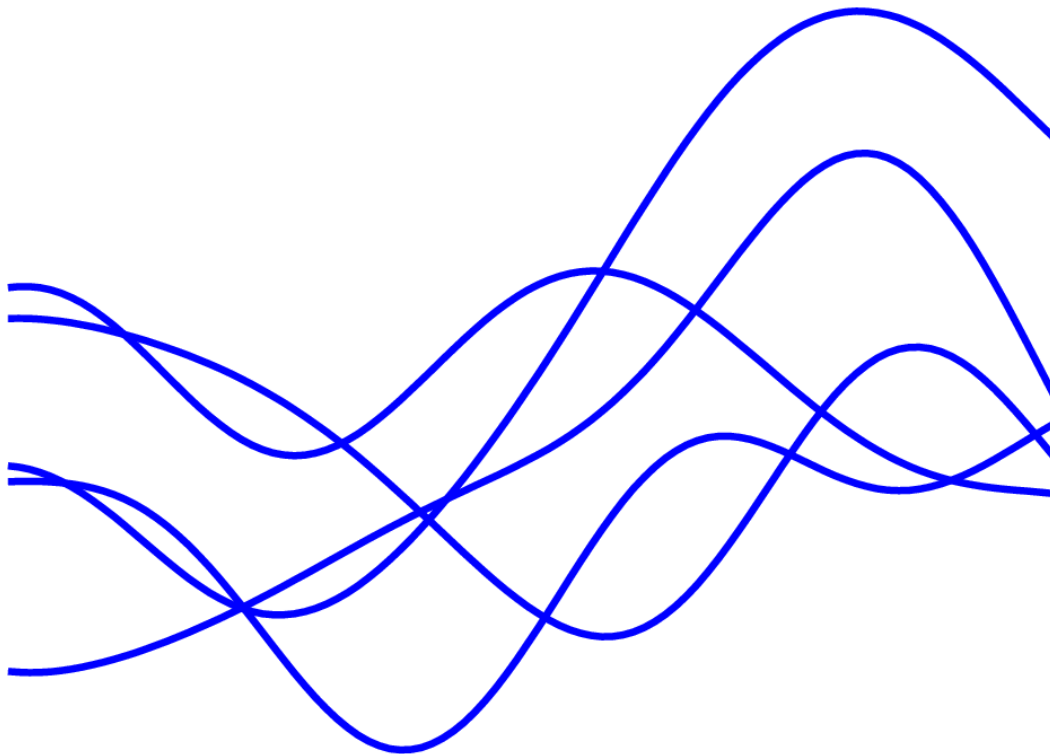
Typical choice in machine learning:  $k(\cdot, \cdot)$  on data from  $\mathbb{R}^D$  uses a parameterized formulation

$$k(\mathbf{x}, \mathbf{x}') = \exp \left( \theta_0 - \sum_{d=1}^D \theta_d (x_d - x'_d)^2 \right)$$

“squared exponential kernel”

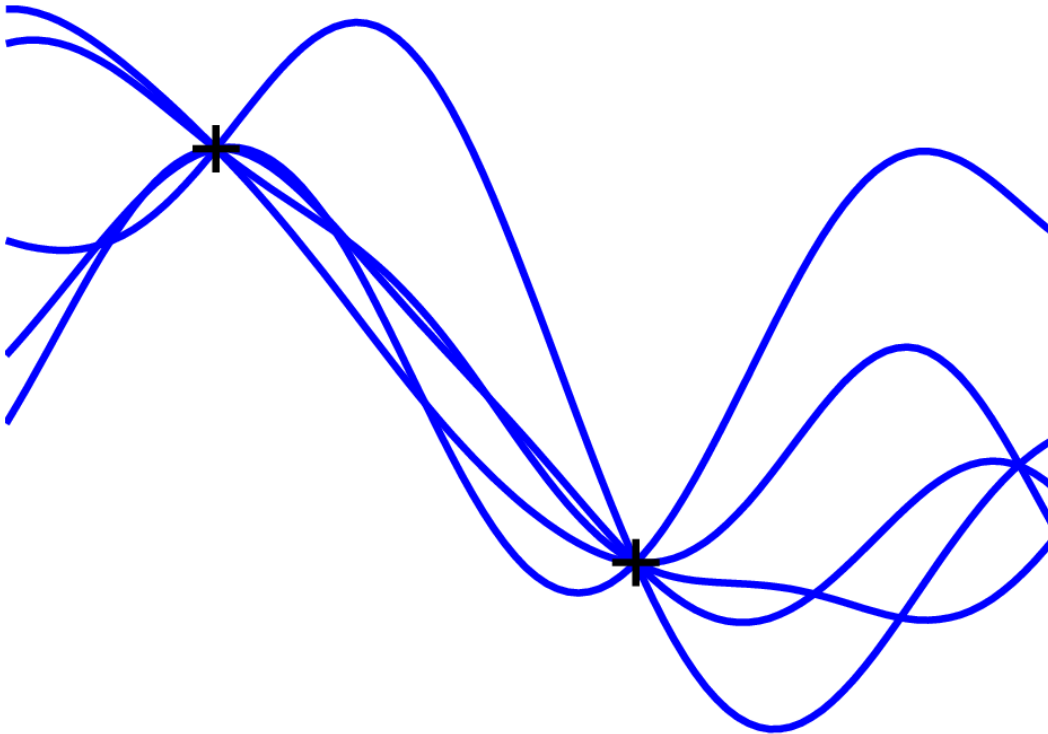
# Gaussian Processes: Samples

5 samples of a GP with squared exponential kernel



# Gaussian Processes: Samples

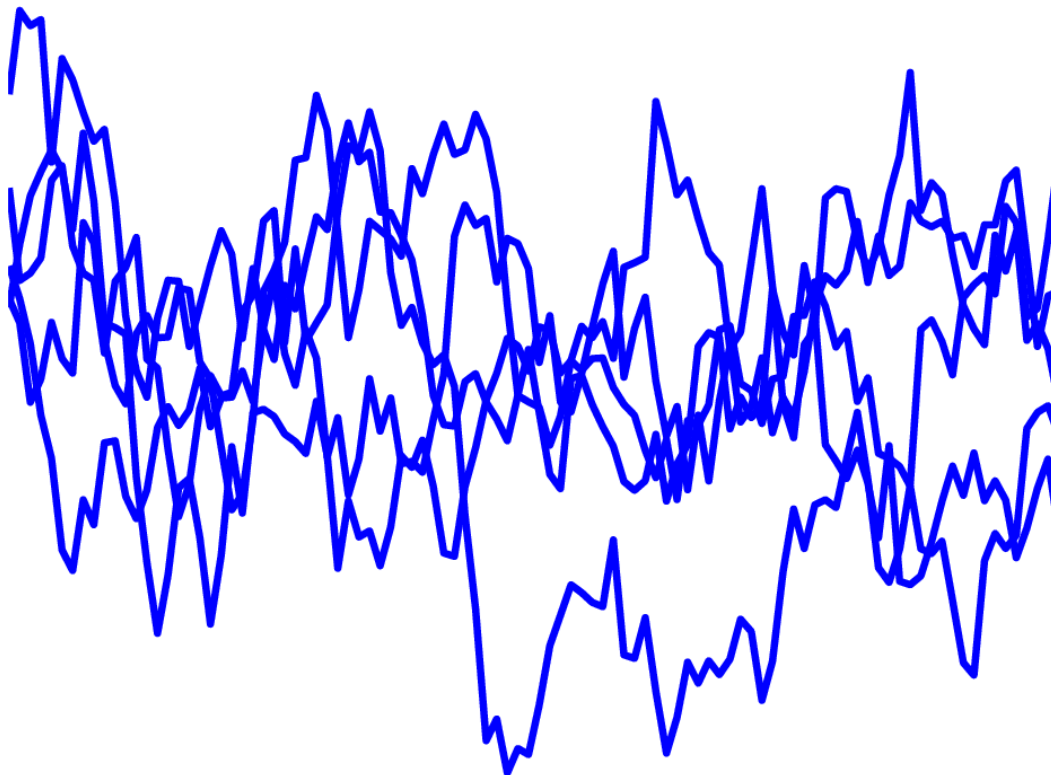
Prediction with Gaussian processes: 5 samples of the GP posterior when two training points are given



# Gaussian Processes: Samples

5 samples of a GP with Ornstein-Uhlenbeck kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-|\mathbf{x} - \mathbf{x}'|)$$





# Bayesian Regression with GPs (1)

We have observed samples from some unknown function  $f$  on locations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , corrupted by Gaussian noise with variance  $\sigma^2$ :

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad i = 1, \dots, N$$

$$\epsilon_i \sim N(0, \sigma^2) \quad \text{Gaussian noise}$$

$$(f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)) \sim N(0, K) \quad \text{Gaussian process prior, } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

$$(y_1, \dots, y_N) \sim N(0, K + \sigma^2 I)$$

We wish to predict at test point  $\mathbf{x}^*$ . From the GP property, we also know that

$$(y_1, \dots, y_N, f(\mathbf{x}^*)) \sim N(0, K^*)$$

Everything is Gaussian: Conditional mean (prediction) is

$$p(f(\mathbf{x}^*) | y_1, \dots, y_N, \theta) = \frac{p(y_1, \dots, y_N, f(\mathbf{x}^*) | \theta)}{p(y_1, \dots, y_N | \theta)}$$

$$E(f(\mathbf{x}^*) | \mathcal{D}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}^*, \mathbf{x}_i)$$

$$\text{with } (K + \sigma^2 \mathbf{1}) \boldsymbol{\alpha} = \mathbf{y}$$

# Evidence maximization

We have so far assumed that the covariance function is given. GPR allows us to infer the most likely parameters of the kernel function from the training data.

Key quantity is the **evidence**: (marginal) log likelihood of the training data under the GP model with kernel hyperparameters  $\theta$ :

$$\log p(\mathbf{y}|\mathbf{x}_1, \dots, \mathbf{x}_N, \theta) = -\frac{1}{2} \log \det(K + \sigma^2 \mathbf{1}) - \frac{1}{2} \mathbf{y}^\top (K + \sigma^2 \mathbf{1})^{-1} \mathbf{y} - \frac{N}{2} \log 2\pi \quad (1)$$

**Integrate out** Assume prior distributions on the hyperparameters  $\theta$  and integrate out with Monte Carlo method

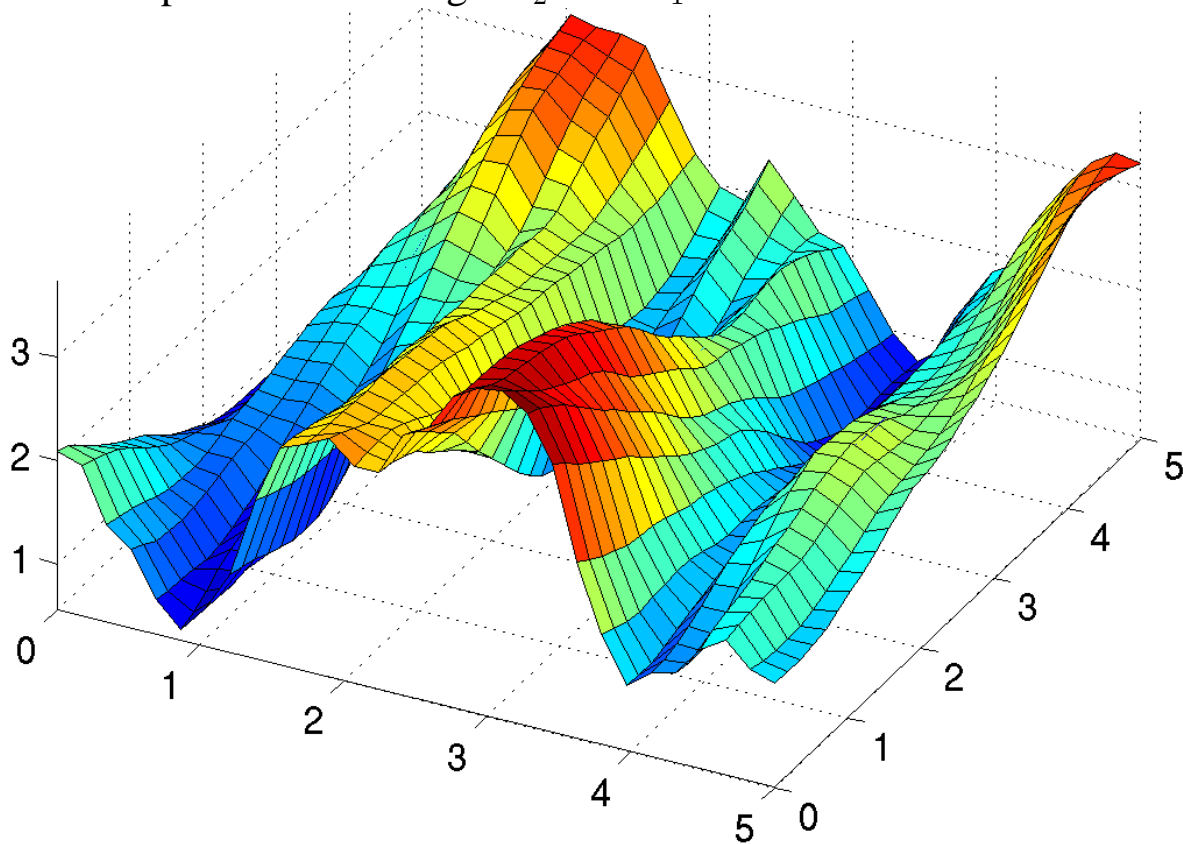
**Maximization** Maximize  $\log p(\mathbf{y}|\mathbf{x}_1, \dots, \mathbf{x}_N, \theta)$  with respect to  $\theta$  (e.g. using conjugate gradient method)

GP with evidence maximization does implicit input pruning: Small  $\theta_d$  in the kernel function indicate irrelevant inputs

$$k(\mathbf{x}, \mathbf{x}') = \exp \left( \theta_0 - \sum_{d=1}^D \theta_d (x_d - x'_d)^2 \right)$$

# Samples of a GP With Anisotropic Kernel

Gaussian process with weight  $\theta_2 = 25\theta_1$ :



## Bayesian Regression with GPs (2)

Cookbook: How to do Gaussian process regression with training data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N, y_1, \dots, y_N\}$

1. Assume a parameterized model for the covariance function  $k_\theta(\cdot, \cdot)$  with parameters  $\theta$
2. Choose parameters  $\theta$  to maximize evidence: Choose the model that best explains the training data
3. Compute kernel matrix  $K$ , with  $K_{ij} = k_\theta(\mathbf{x}_i, \mathbf{x}_j)$
4. Solve linear system  $(K + \sigma^2 \mathbf{1})\alpha = \mathbf{y}$
5. Predict value  $y^*$  for a test point  $\mathbf{x}^*$  from  $\sum_{i=1}^N \alpha_i k_\theta(\mathbf{x}^*, \mathbf{x}_i)$  (mean of the GP posterior)
6. Obtain error bars for  $y^*$  from the variance of the GP posterior

# Relations to Other Methods

---

Regression with Gaussian Processes is closely related to

**Kriging** Best linear unbiased estimator (minimum variance), originally developed for spatial data in geostatistics

**Infinite Neural Networks** A 2-layer neural network with increasing number of hidden units converges to a Gaussian process predictor with a certain type of covariance function

**Regularization** Choice of kernel corresponds to a regularization operator (smoothness prior)

# Overview

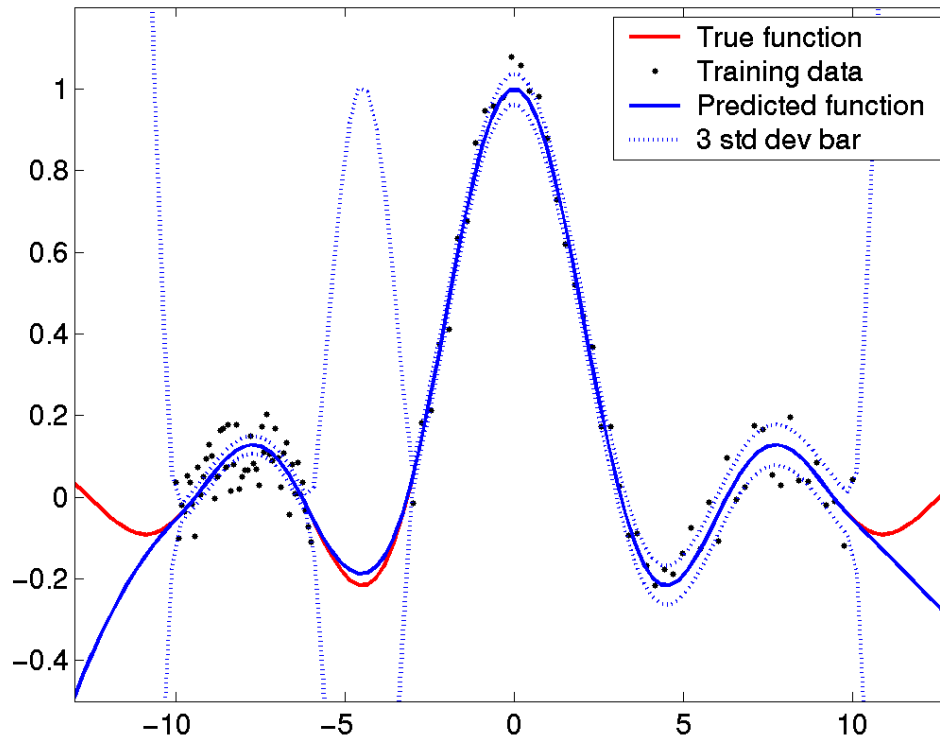
---

- Gaussian processes basics
- Gaussian processes for Bayesian regression (GPR)
- The evidence procedure in GPR
- **Choice of kernel functions**
- Application areas

# Gaussian Processes versus RVM

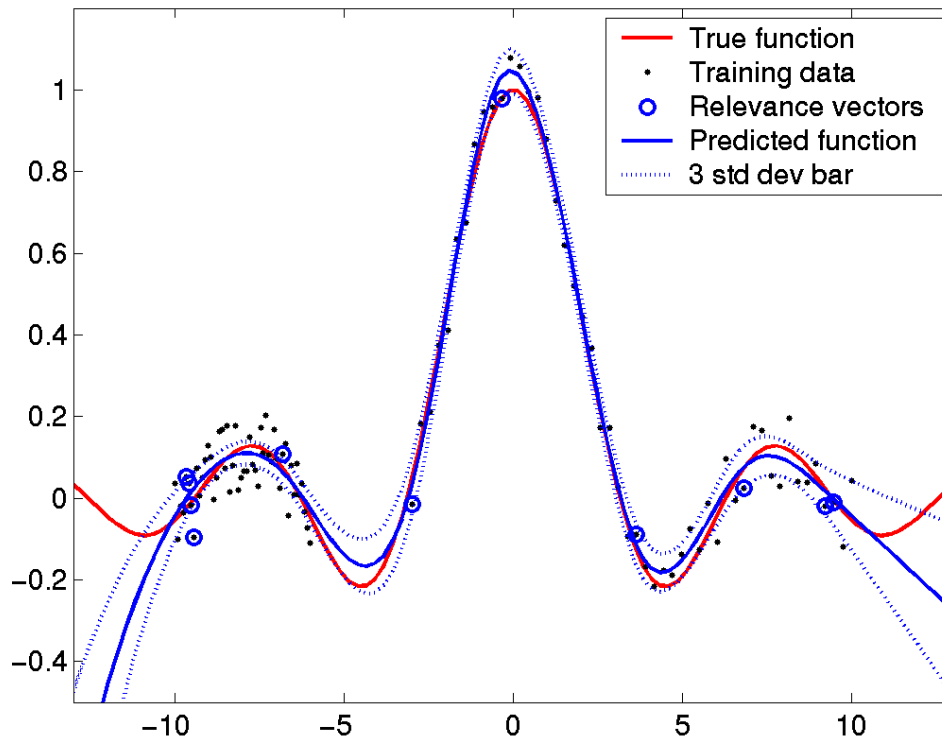
- Noisy sinc problem with regions of missing data
- Thin plate spline kernel function (as in Tipping's toy example)

True function and GP prediction with error bars:



# Gaussian Processes versus RVM

True function and RVM prediction with error bars:



RVM is grossly over-optimistic about the variance!



# On the Choice of Kernel Function (1)

Widely used kernel function: Squared exponential (aka Gaussian, aka RBF kernel)

$$k(\mathbf{x}, \mathbf{x}') = \exp \left( \theta_0 - \sum_{d=1}^D \theta_d (x_d - x'_d)^2 \right) = \exp \left( \theta_0 - \|\mathbf{x} - \mathbf{x}'\|^2 \right)$$

Criticism:

- Unreasonable smoothness assumption
- Variance of predictions tends to be underestimated

Alternatives: **Matern kernel**

$$k(\mathbf{x}, \mathbf{x}') = M_\nu(z) = \frac{2(\sqrt{\nu}z)^\nu}{\Gamma(\nu)} K_\nu(2\sqrt{\nu}z)$$

with  $z = \|\mathbf{x} - \mathbf{x}'\|$ .  $\Gamma(\nu)$  is the Gamma function and  $K_\nu(r)$  is the modified Bessel function of the second kind of degree  $\nu$ .

## On the Choice of Kernel Function (2)

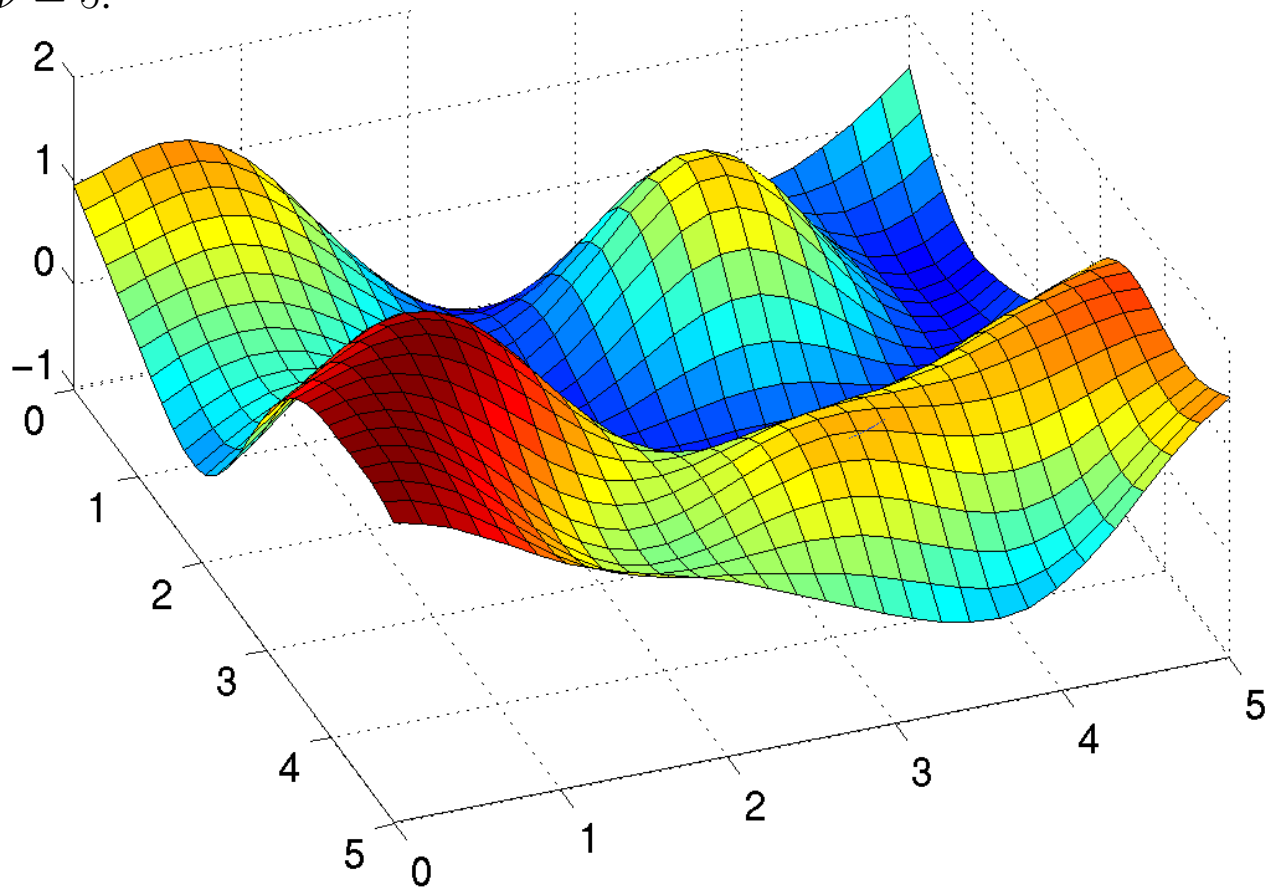
Matern kernel

$$k(\mathbf{x}, \mathbf{x}') = M_\nu(z) = \frac{2(\sqrt{\nu}z)^\nu}{\Gamma(\nu)} K_\nu(2\sqrt{\nu}z)$$

- With  $\nu \rightarrow \infty$ , the Matern kernel converges to the squared exponential kernel
- $\nu$  allows a continuous parameterization of the **fractal dimension** of the underlying process from smooth ( $\nu \rightarrow \infty$ ) to rough
- Often, there is no basis for knowing the degree of smoothness of some process *a priori*. The Matern kernel allows to **infer the smoothness from the data** through the evidence procedure.
- Disadvantage: Computationally more intensive (Bessel and Gamma functions), numerically sensitive

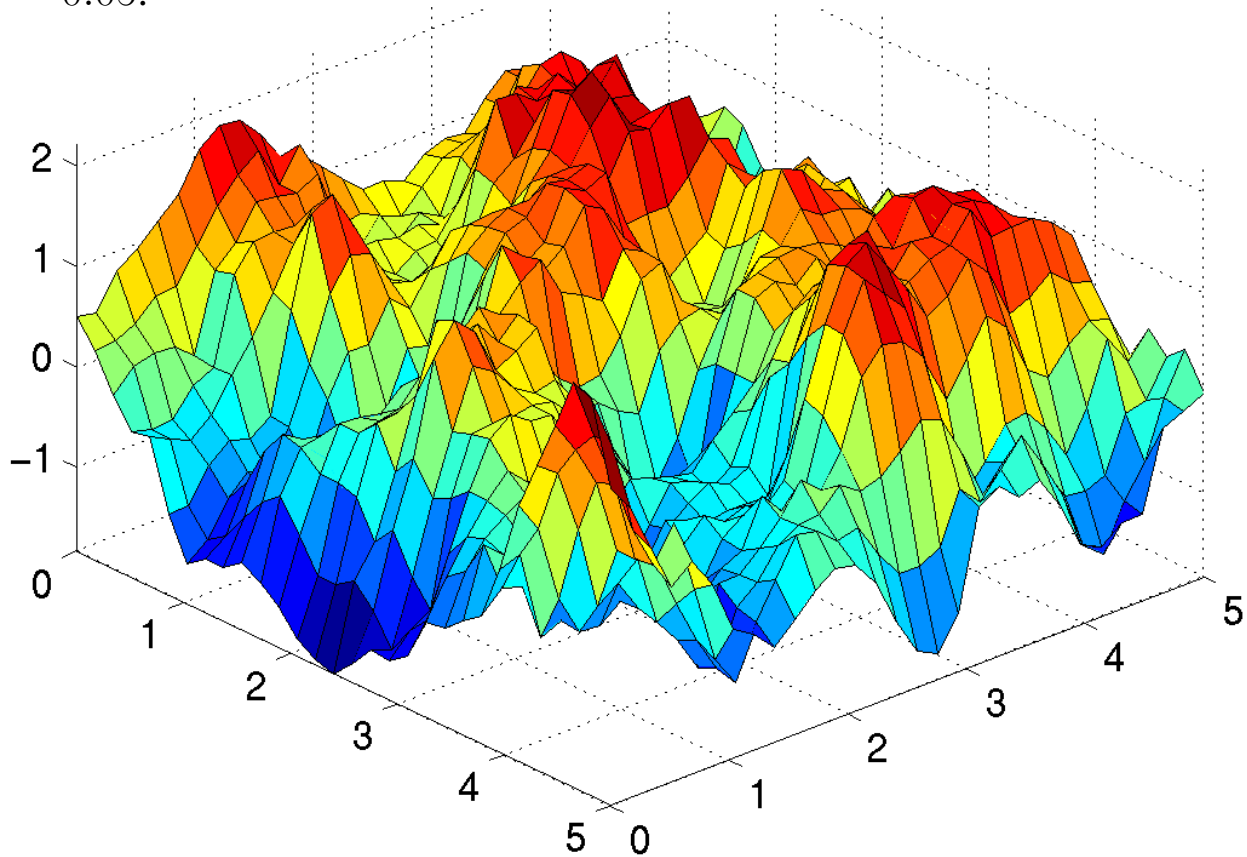
# GP with Matern Kernel: Samples (1)

$\nu = 5$ :



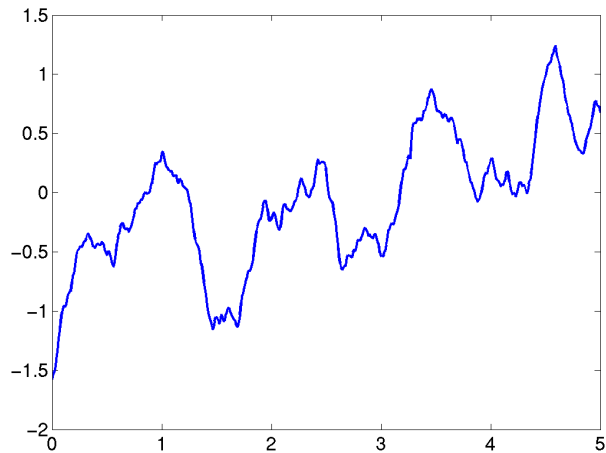
# GP with Matern Kernel: Samples (2)

$\nu = 0.05$ :

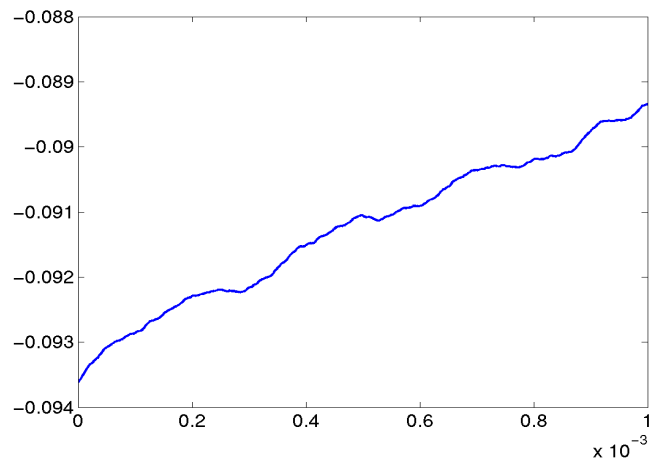


# GP with Matern Kernel: Samples (3)

Normal scale



Zoom in



- Sample paths become increasingly rough when  $\nu \rightarrow 0$
- Random structure visible at all scales

# Application areas (1)

- Regression and classification problems
- Modelling of dynamic systems
- Neural responses (finding the most likely stimulus to a pool of spiking neurons)
- Learning kernel functions (Platt, MSR: Kernels for music playlists)

## Advantages:

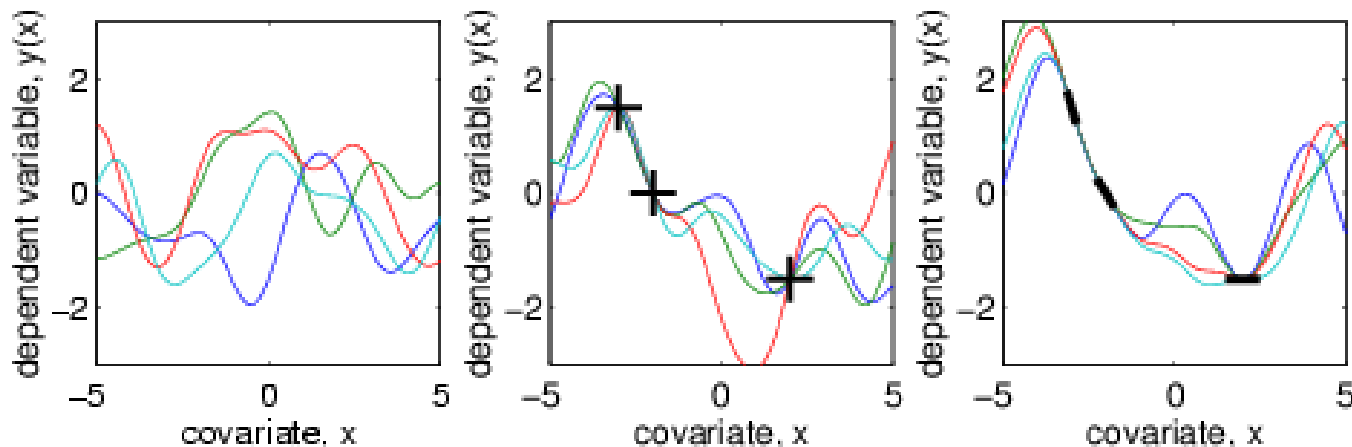
- Tractable exact Bayesian inference: Many of the distributions of interest are again Gaussian
- Kernel functions allow working on arbitrary data types (as long as you can define a kernel on them)

# Application Areas (2): Dynamic Systems

- Derivative of a Gaussian process is again a Gaussian process:

$$\text{cov} \left( f(\mathbf{x}), \frac{\partial f(\mathbf{x}')}{\partial x_i} \right) = \frac{\partial}{\partial x_i} \text{cov}(f(\mathbf{x}), f(\mathbf{x}'))$$

- Derivative observations can be incorporated into a GP model in a straightforward way, leading to reduced predictive variance
- Application: Build an accurate model of complex probability distributions from few samples, for increasing speed of Monte Carlo simulations (Rasmussen, Bayesian Statistics vol. 7)



## Application Areas (3)

### Regression problems on non-vectorial data

- Gaussian processes only require knowledge of kernel evaluations  $k(\mathbf{x}, \mathbf{x}')$
- We can work with any type of data we can formulate a kernel for (similar to Support Vector Machines)
- Application: Gaussian processes for generating music playlists from discrete features (Platt, NIPS\*2001)



# Conclusions

---

- Gaussian processes are an elegant and easy-to-use way for doing nonlinear regression
- GPs have shown excellent performance on many regression and classification tasks
- General feature of kernel systems: handling of arbitrary data types (as long as you can define a kernel = covariance function for them)

*Active areas of research:*

- Kernel design, in particular for applications with non-vectorial data