

# ENM 531: Data-driven Modeling and Probabilistic Scientific Computing

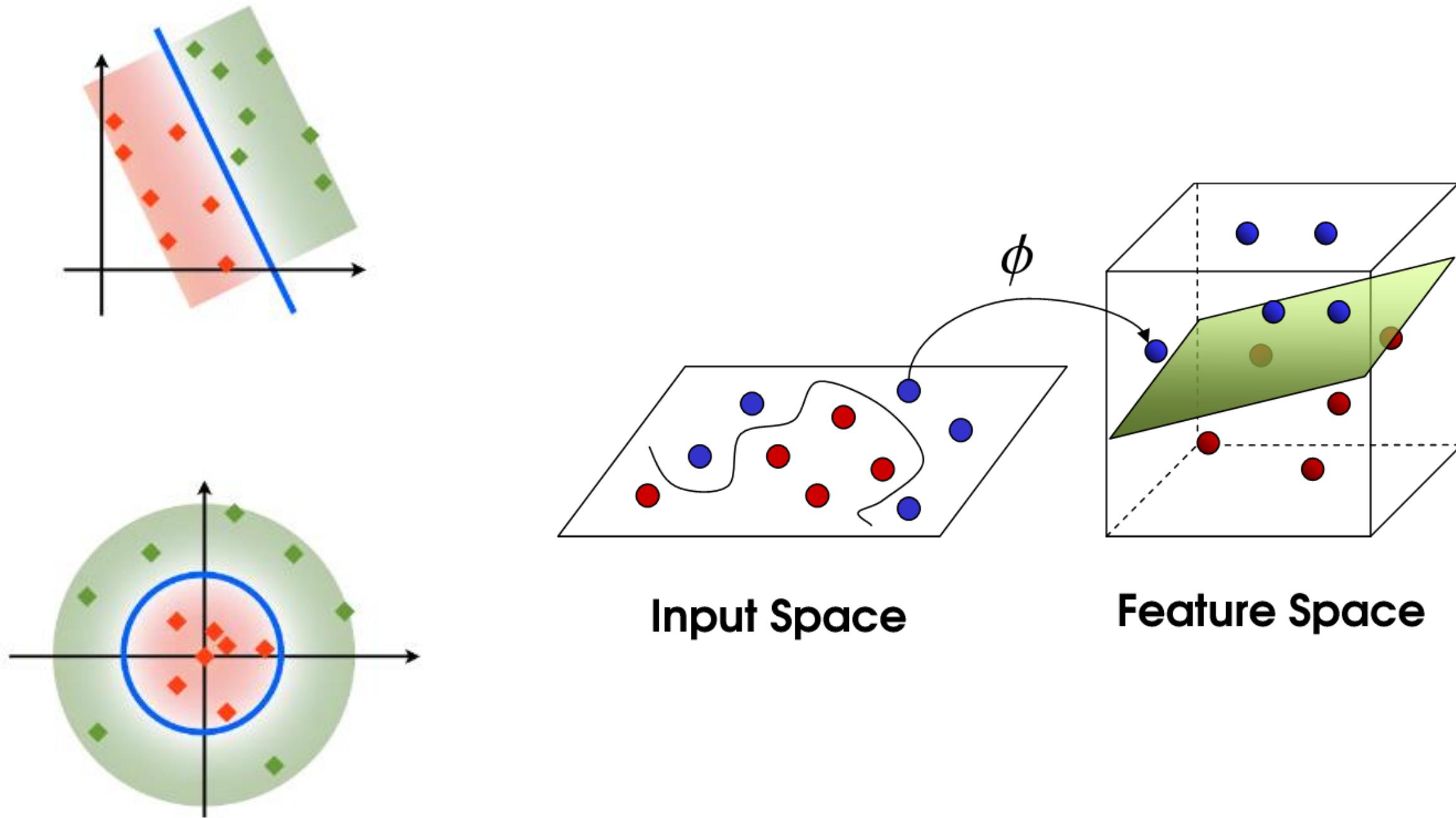
## *Lecture #18: Gaussian process regression*

April 6, 2023

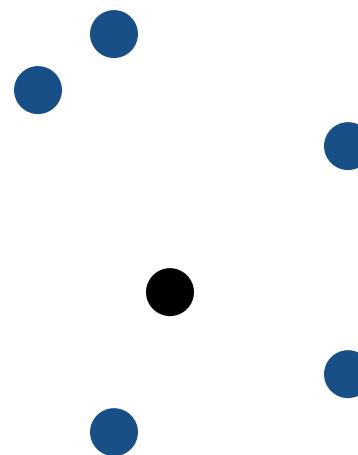


Recall: “Linearization” by embedding to higher dimensions

$$f(\mathbf{x}) = \langle \theta, \phi(\mathbf{x}) \rangle_{\mathcal{H}}, \quad \phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$



# Kernel methods



$$f(\mathbf{x}) = \langle \theta, \phi(\mathbf{x}) \rangle_{\mathcal{H}}, \quad \phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

$$m \rightarrow \infty$$

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$$

$$f(\mathbf{x}) = \sum_{i=1}^n (\mathbf{K}^{-1} \mathbf{y})_i k(\mathbf{x}_i, \mathbf{x}), \quad \mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

Key take-away points:

- Kernels and representer theorems: learning with infinite-dimensional linear models can be done in time that depend on the number of observations by using a kernel function.
- Kernels on  $\mathbb{R}^d$ : such models include polynomials and classical Sobolev spaces (functions with square-integrable partial derivatives).
- Algorithms: convex optimization algorithms can be applied with theoretical guarantees and many dedicated developments to avoid the quadratic complexity of computing the kernel matrix.
- Analysis of well-specified models: When the target function is in the associated function space, learning can be done with rates that are independent of dimension.
- Analysis of mis-specified models: if the target is not in the RKHS, the curse of dimensionality cannot be avoided in the worst case situations of few existing derivatives of the target function, but the methods are adaptive to any amount of intermediate smoothness.
- Sharp analysis of ridge regression: for the square loss, a more involved analysis leads to optimal rates in a variety of situations in  $\mathbb{R}^d$ .

# Kernel methods

$$f(\mathbf{x}) = \langle \theta, \phi(\mathbf{x}) \rangle_{\mathcal{H}}, \quad \phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

$$m \rightarrow \infty$$



$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$$

$$f(\mathbf{x}) = \sum_{i=1}^n (\mathbf{K}^{-1} \mathbf{y})_i k(\mathbf{x}_i, \mathbf{x}), \quad \mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$



$$p(f|\mathcal{D}) = \frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})}$$

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$$

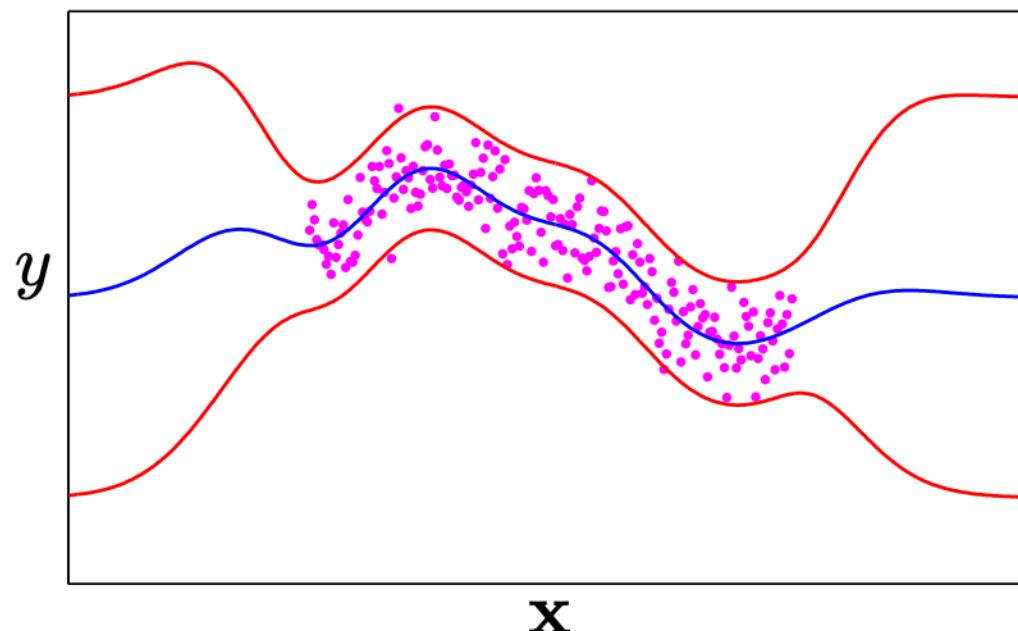
Key take-away points:

- Kernels and representer theorems: learning with infinite-dimensional linear models can be done in time that depend on the number of observations by using a kernel function.
- Kernels on  $\mathbb{R}^d$ : such models include polynomials and classical Sobolev spaces (functions with square-integrable partial derivatives).
- Algorithms: convex optimization algorithms can be applied with theoretical guarantees and many dedicated developments to avoid the quadratic complexity of computing the kernel matrix.
- Analysis of well-specified models: When the target function is in the associated function space, learning can be done with rates that are independent of dimension.
- Analysis of mis-specified models: if the target is not in the RKHS, the curse of dimensionality cannot be avoided in the worst case situations of few existing derivatives of the target function, but the methods are adaptive to any amount of intermediate smoothness.
- Sharp analysis of ridge regression: for the square loss, a more involved analysis leads to optimal rates in a variety of situations in  $\mathbb{R}^d$ .

# Nonlinear regression

Consider the problem of **nonlinear regression**:

You want to learn a **function  $f$**  with **error bars** from **data  $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$**



A **Gaussian process** defines a distribution over functions  $p(f)$  which can be used for Bayesian regression:

$$p(f|\mathcal{D}) = \frac{p(f)p(\mathcal{D}|f)}{p(\mathcal{D})}$$

# Carl Friedrich Gauss (1777–1855)

Paying Tolls with A Bell

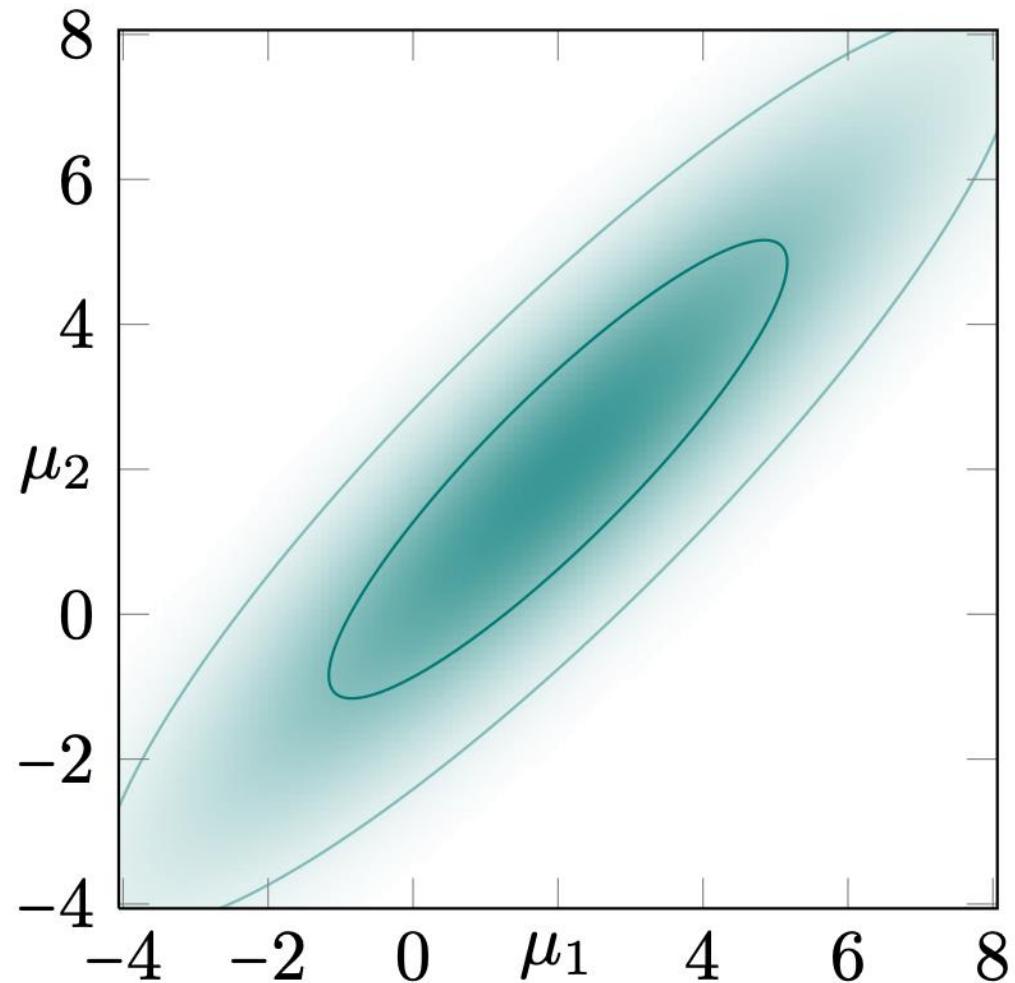
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



# The Gaussian distribution

## Multivariate Form

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right]$$

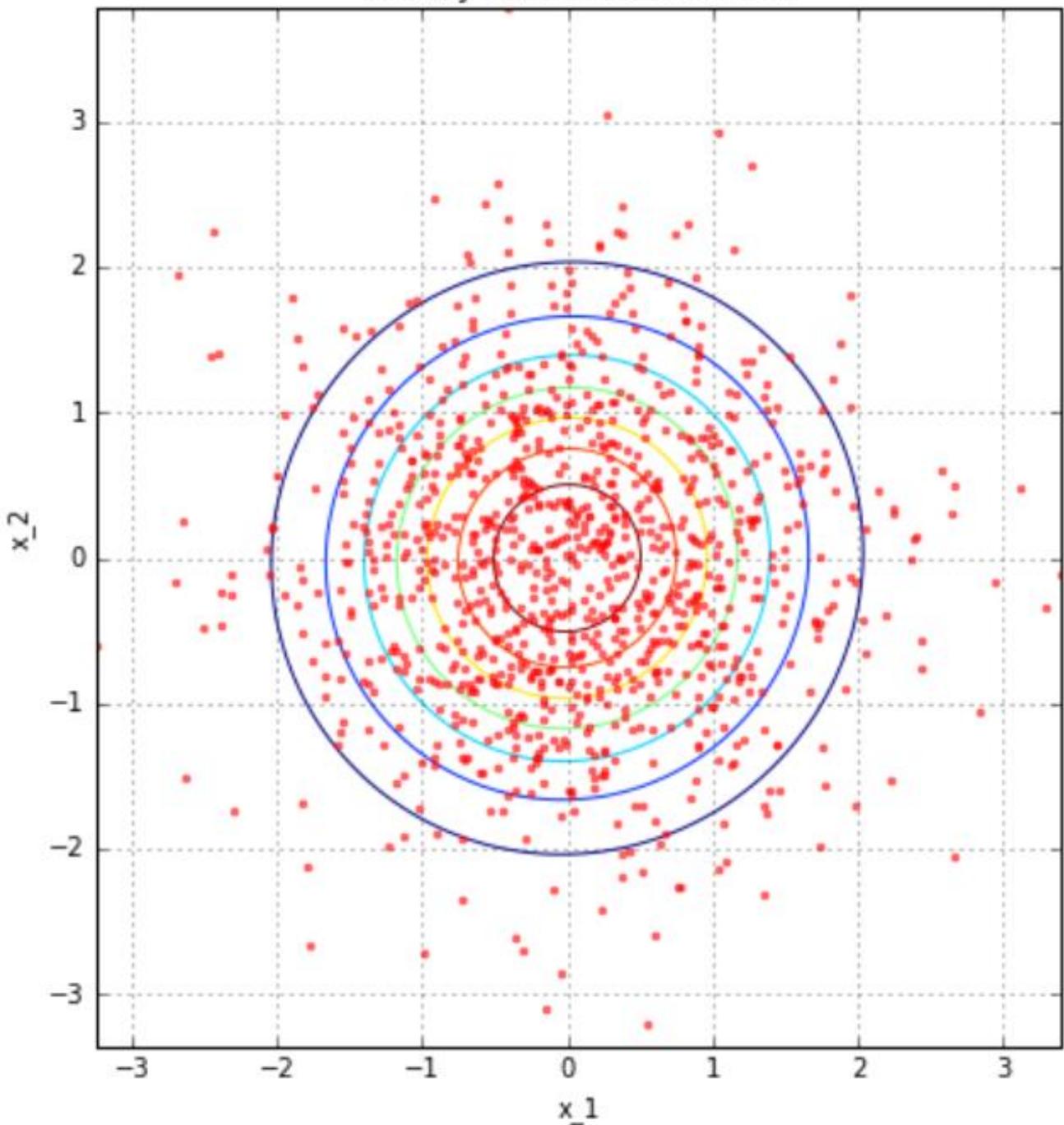


- ▶  $x, \mu \in \mathbb{R}^N, \Sigma \in \mathbb{R}^{N \times N}$
- ▶  **$\Sigma$  is positive semidefinite, i.e.**
  - ▶  $v^\top \Sigma v \geq 0$  for all  $v \in \mathbb{R}^N$
  - ▶ Hermitian, all eigenvalues  $\geq 0$

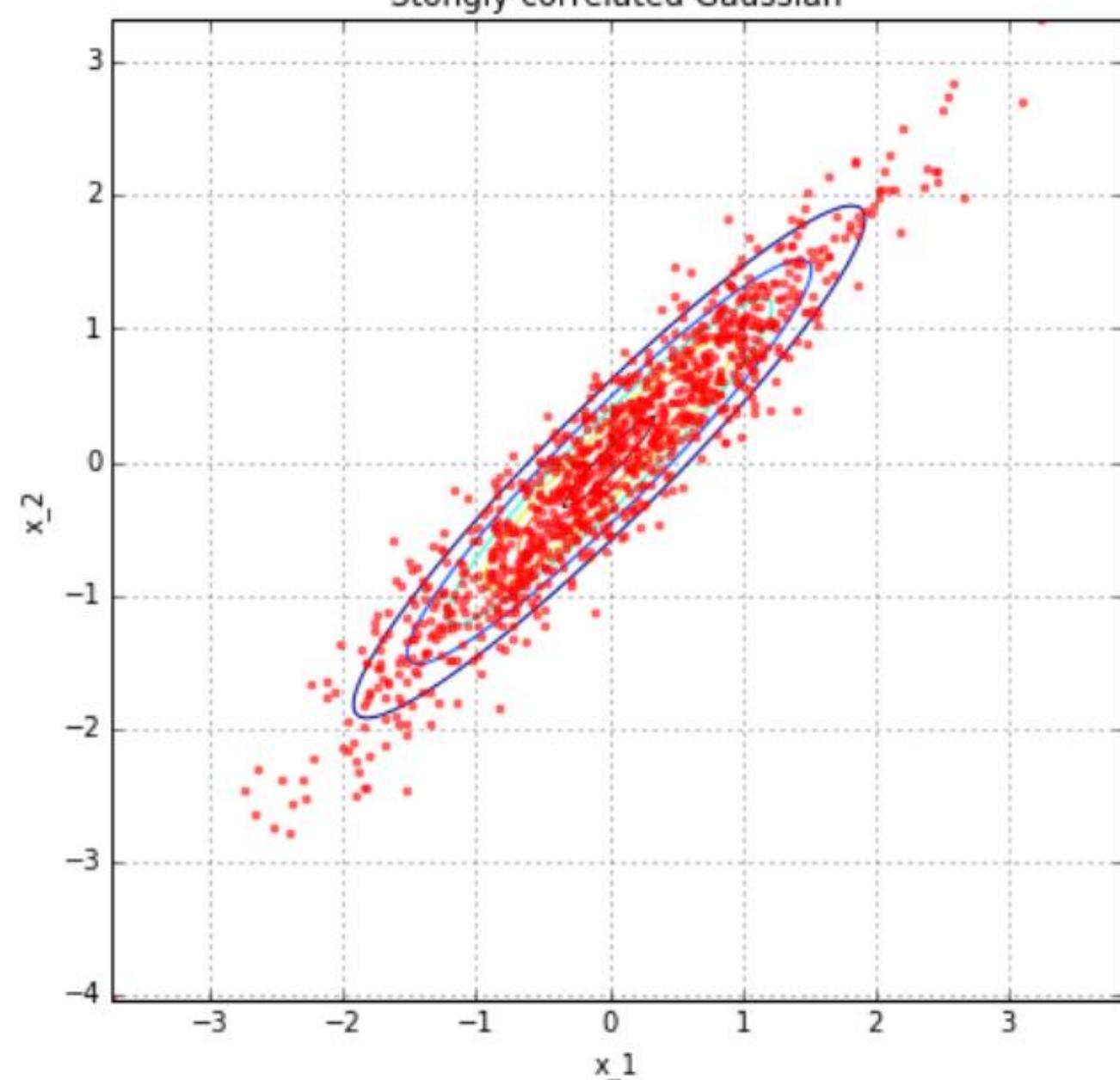
DEMO: <https://distill.pub/2019/visual-exploration-gaussian-processes/>

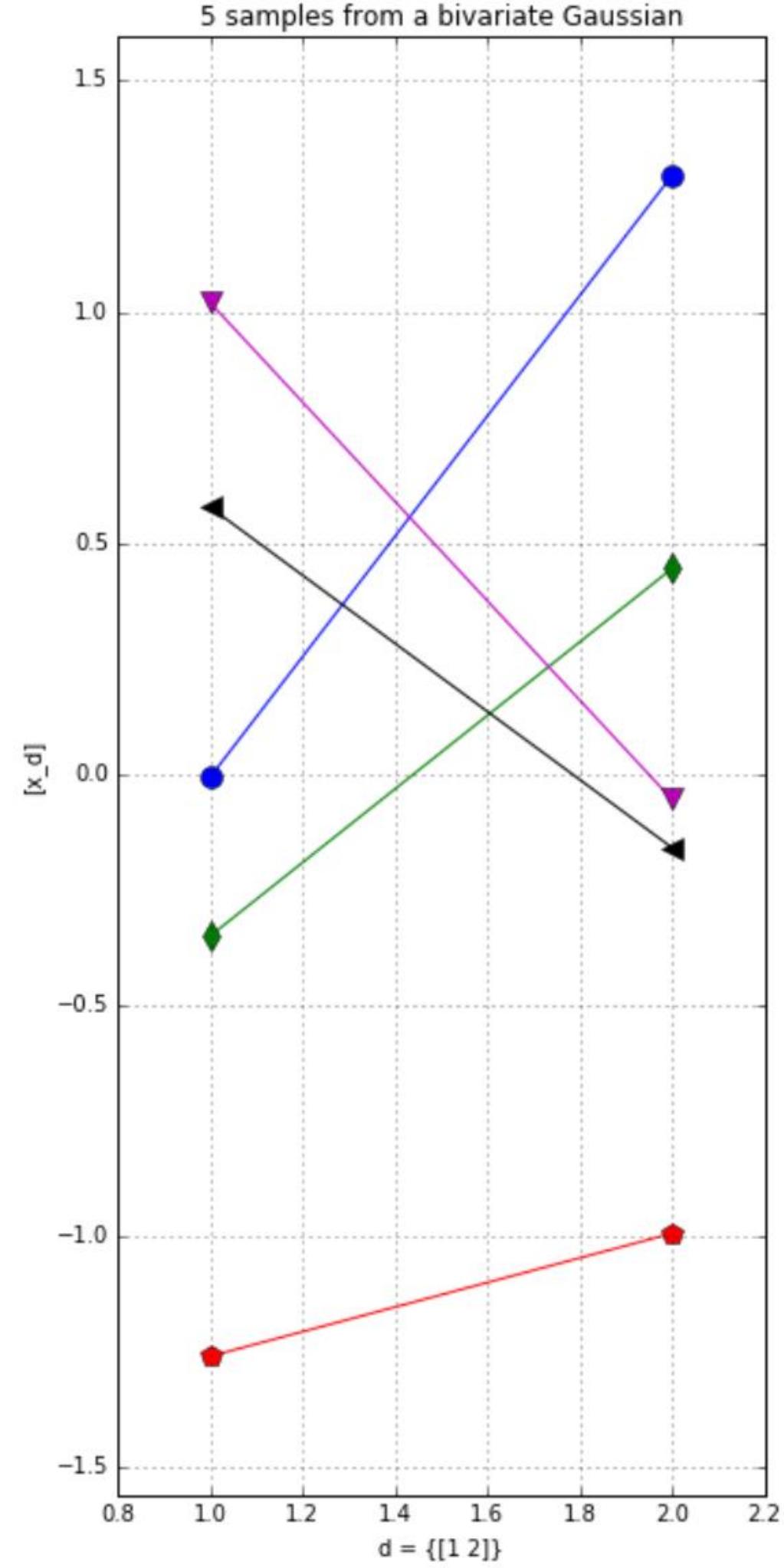
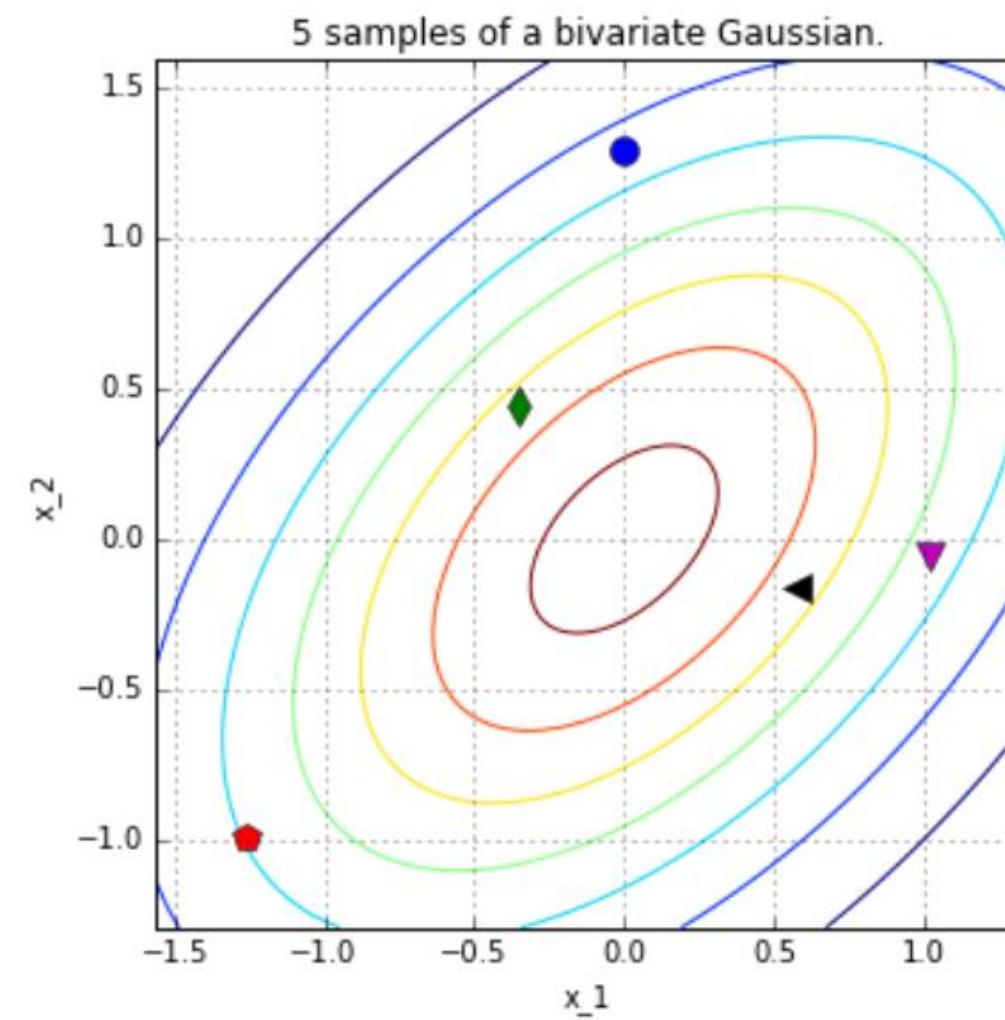
[http://mlss.tuebingen.mpg.de/2013/2013/hennig\\_slides1.pdf](http://mlss.tuebingen.mpg.de/2013/2013/hennig_slides1.pdf)

Weakly correlated Gaussian

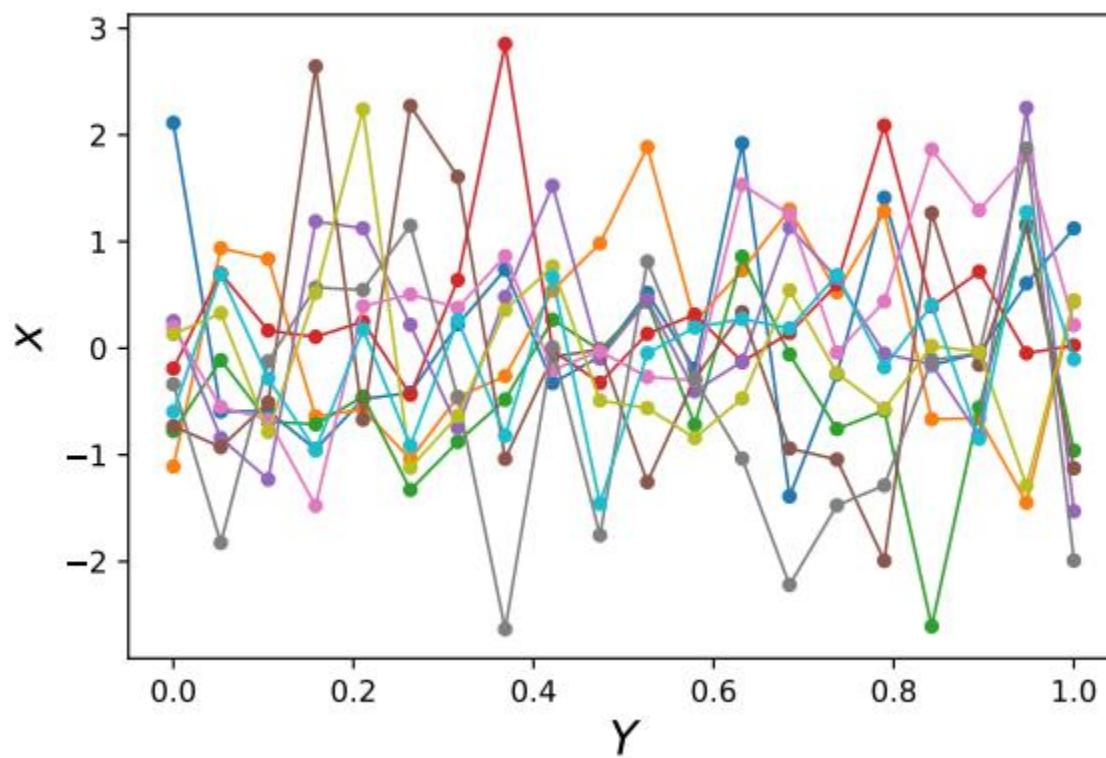


Strongly correlated Gaussian



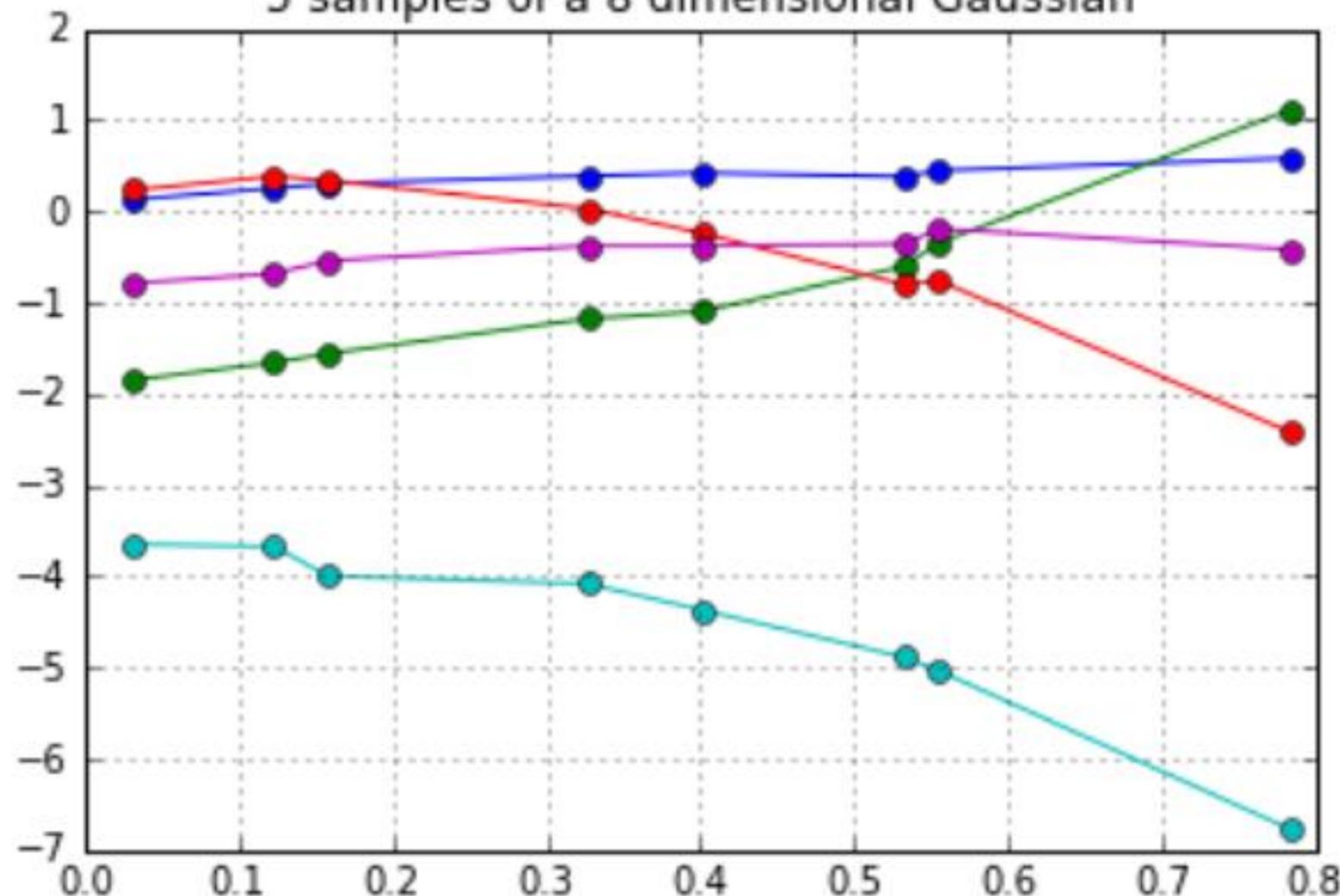


# Draws from a 20-dim Distribution

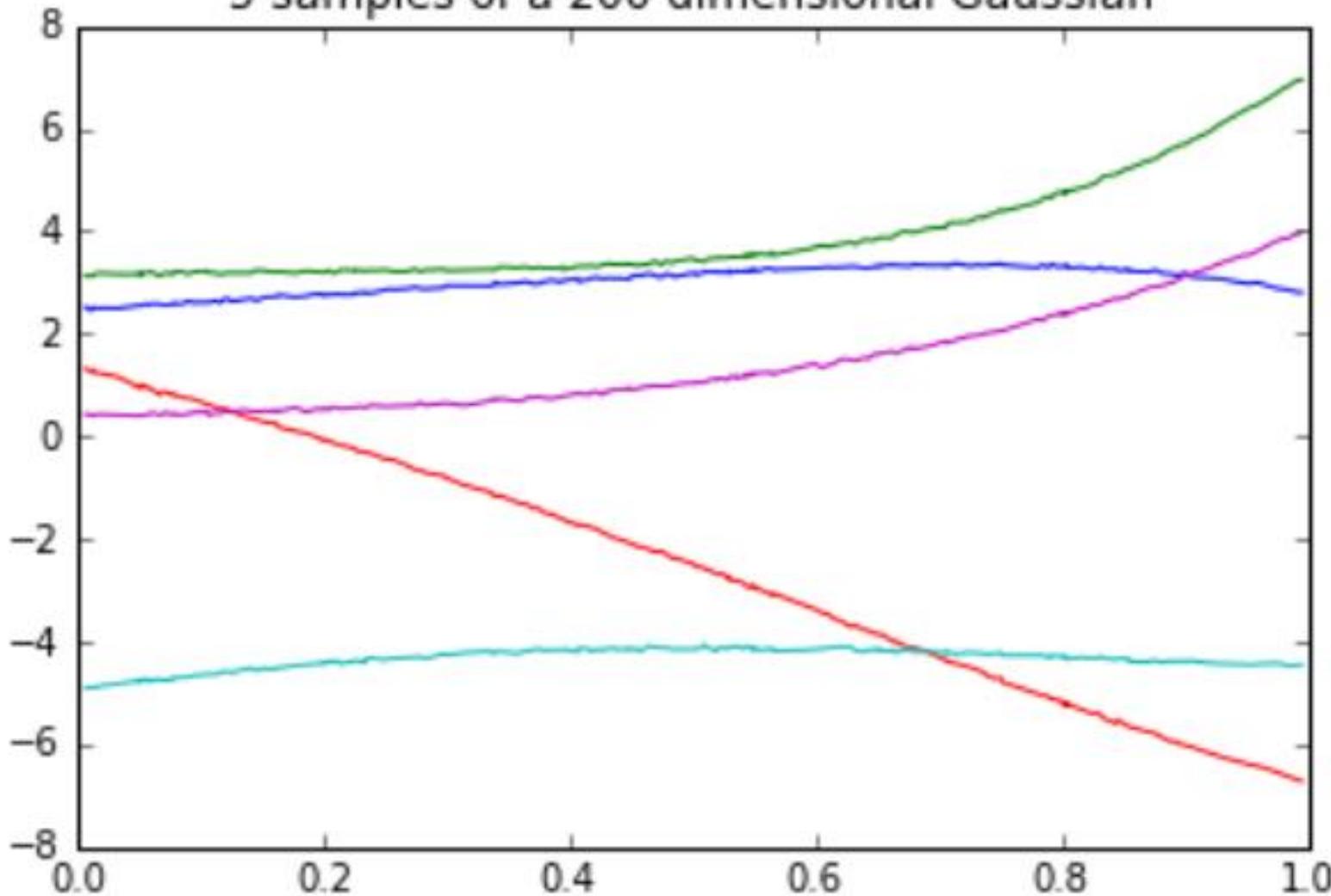


Need indices that are closer together to be correlated.

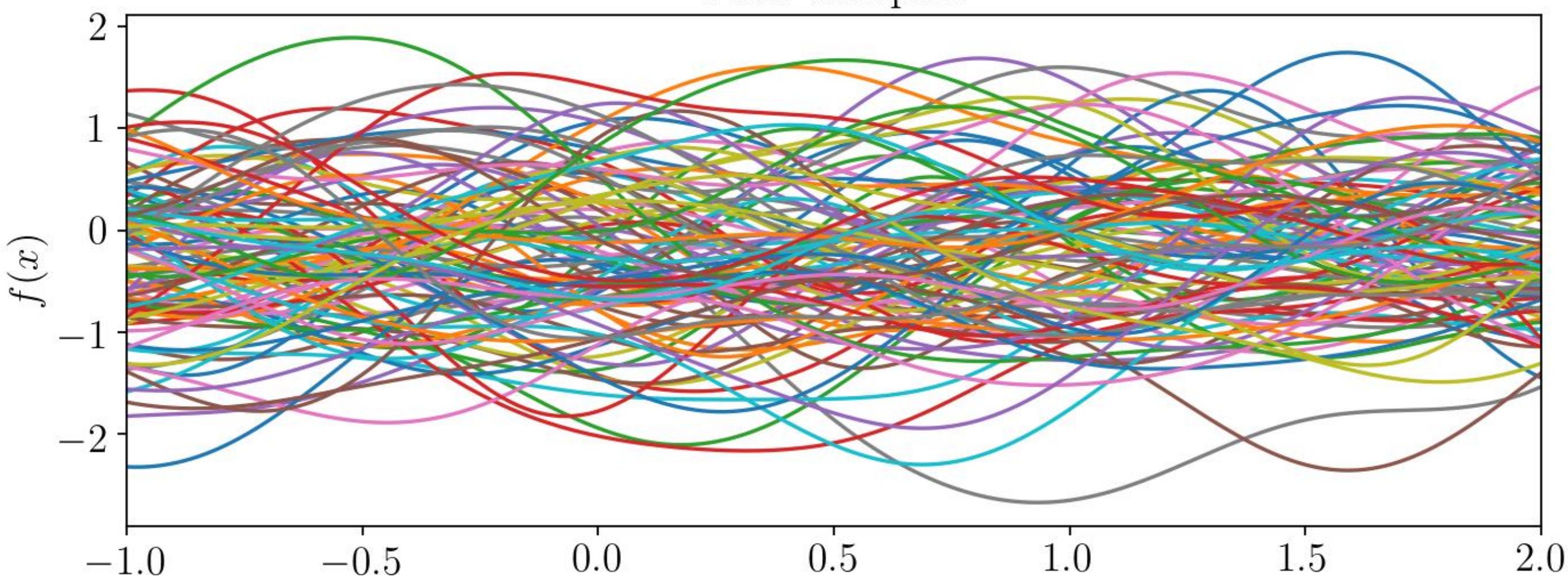
5 samples of a 8 dimensional Gaussian



5 samples of a 200 dimensional Gaussian



Prior samples



## From linear regression to GPs:

- Linear regression with inputs  $x_i$  and outputs  $y_i$ : 
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
- Linear regression with  $M$  basis functions: 
$$y_i = \sum_{m=1}^M \beta_m \phi_m(x_i) + \epsilon_i$$
- Bayesian linear regression with basis functions:  
$$\beta_m \sim \mathcal{N}(\cdot | 0, \lambda_m) \quad (\text{independent of } \beta_\ell, \forall \ell \neq m), \quad \epsilon_i \sim \mathcal{N}(\cdot | 0, \sigma^2)$$
- Integrating out the coefficients,  $\beta_j$ , we find:

$$E[y_i] = 0, \quad Cov(y_i, y_j) = K_{ij} \stackrel{\text{def}}{=} \sum_{m=1}^M \lambda_m \phi_m(x_i) \phi_m(x_j) + \delta_{ij} \sigma^2$$

This is a Gaussian process with covariance function  $K(x_i, x_j) = K_{ij}$ .

This GP has a finite number ( $M$ ) of basis functions. Many useful GP kernels correspond to infinitely many basis functions (i.e. infinite-dim feature spaces).

A multilayer perceptron (neural network) with infinitely many hidden units and Gaussian priors on the weights  $\rightarrow$  a GP (Neal, 1996)

To minimize the loss function

$$L(w) = \sum_{i=1}^N (w^T \phi(x_i) - y) + \frac{\lambda}{2} w^T w$$

given basis vectors  $\{\phi_i\}_{i=1}^M$ , we had that the coefficients  $w$  were given by

$$w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

### Lemma

For any matrix  $A \in \mathbb{R}^{N \times d}$ ,  $y \in \mathbb{R}^d$ , and  $\lambda > 0$ ,

$$(A^T A + \lambda I)^{-1} A^T y = A^T (A A^T + \lambda I)^{-1} y. \quad (2)$$

- The left hand side is solution to normal equation, which means  $A^T A \theta + \lambda \theta = A^T y$ .
- Rearrange terms gives  $\theta = A^T [\frac{1}{\lambda}(y - A\theta)]$ .
- Define  $\alpha = \frac{1}{\lambda}(y - A\theta)$ , then  $\theta = A^T \alpha$ .
- Substitute  $\theta = A^T \alpha$  into  $\alpha = \frac{1}{\lambda}(y - A\theta)$ , we have

$$\alpha = \frac{1}{\lambda}(y - A A^T \alpha).$$

- Rearrange terms gives  $(A A^T + \lambda I)\alpha = y$ , which yields  $\alpha = (A A^T + \lambda I)^{-1} y$ .
- Substitute into  $\theta = A^T \alpha$  gives  $\theta = A^T (A A^T + \lambda I)^{-1} y$ .

## Lemma

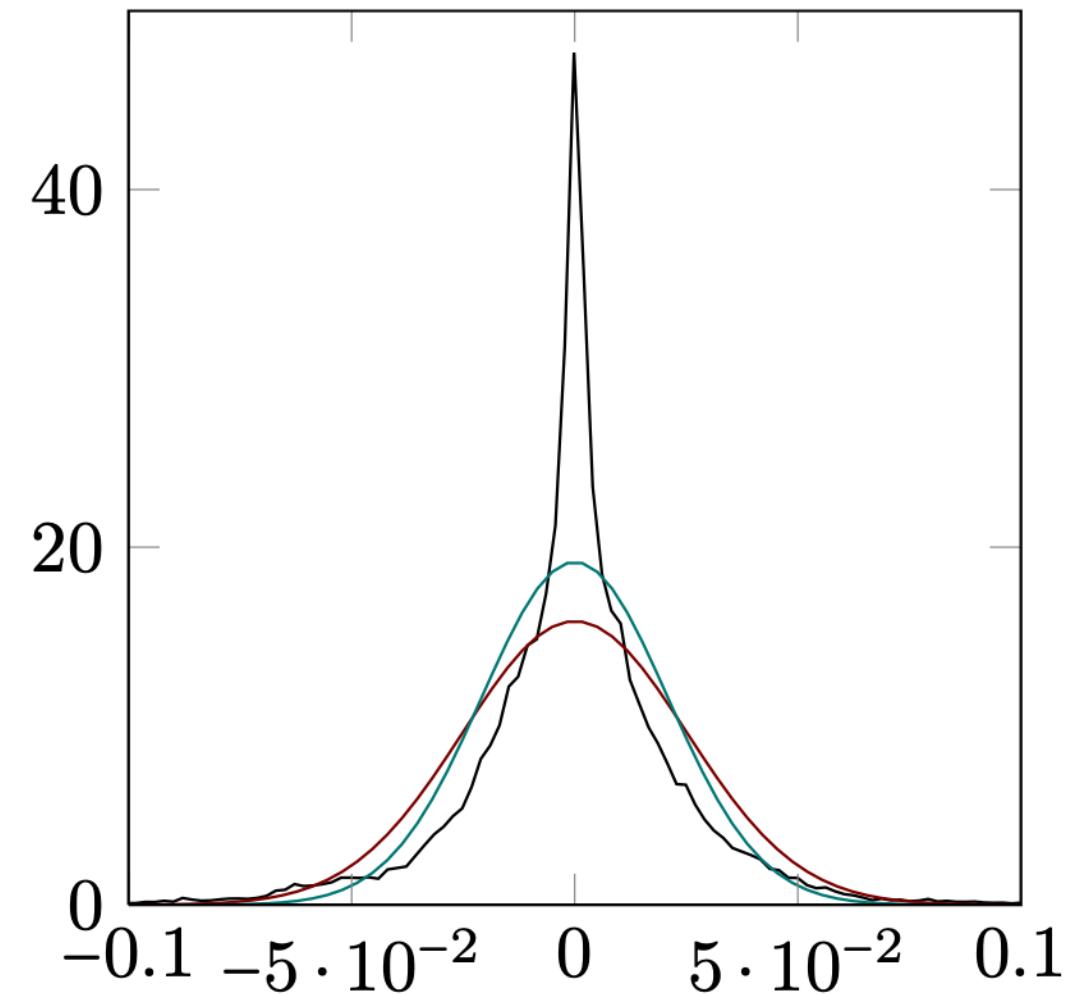
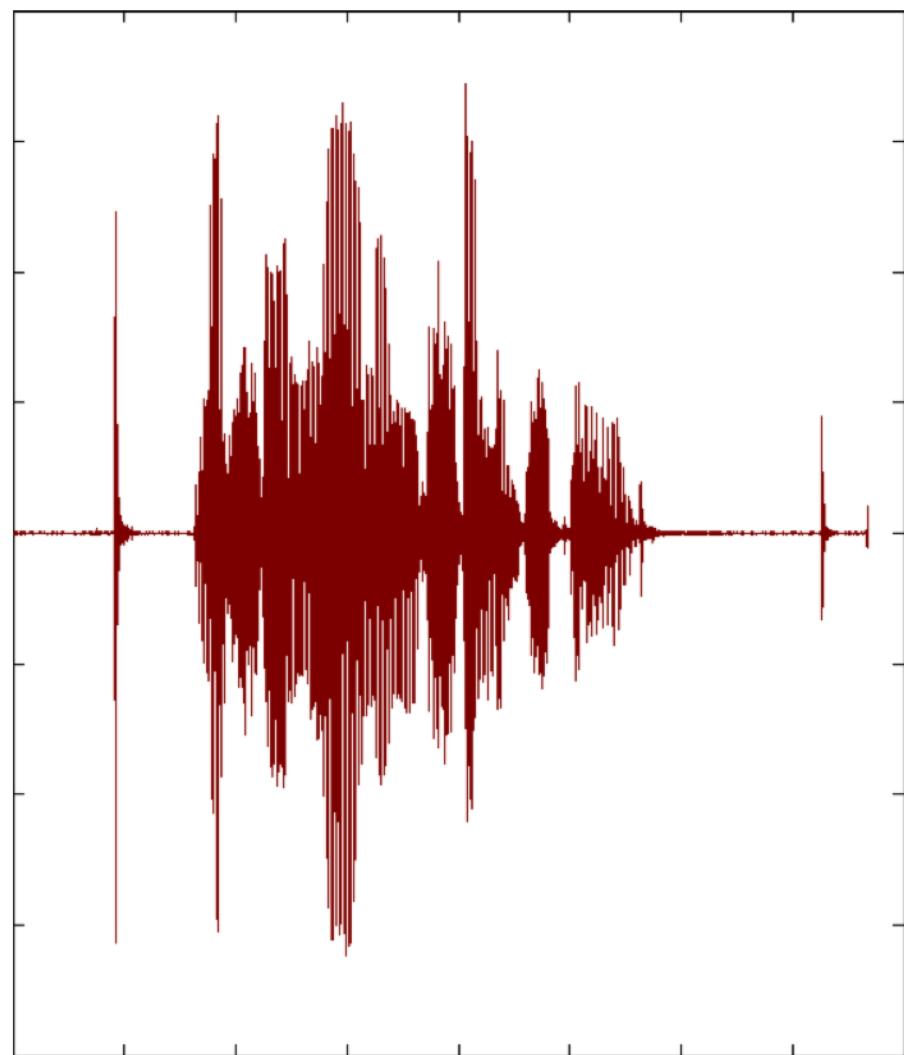
For any matrix  $\mathbf{A} \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^d$ , and  $\lambda > 0$ ,

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}.$$

- Note that in the *primal* formulation we are inverting an  $M \times M$  matrix.
- However, we're inverting an  $N \times N$  matrix in the case of the *dual* problem.  
(typically  $N \gg M$ )
- However, the advantage of the dual formulation is that we can directly work with kernels.
- Avoids the explicit introduction of feature vectors which allows us to use feature spaces of high, even infinite dimensionality.

# Why Gaussian?

an experiment



- ▶ nothing in the real world is Gaussian (except sums of i.i.d. variables)
- ▶ But nothing in the real world is **linear** either!

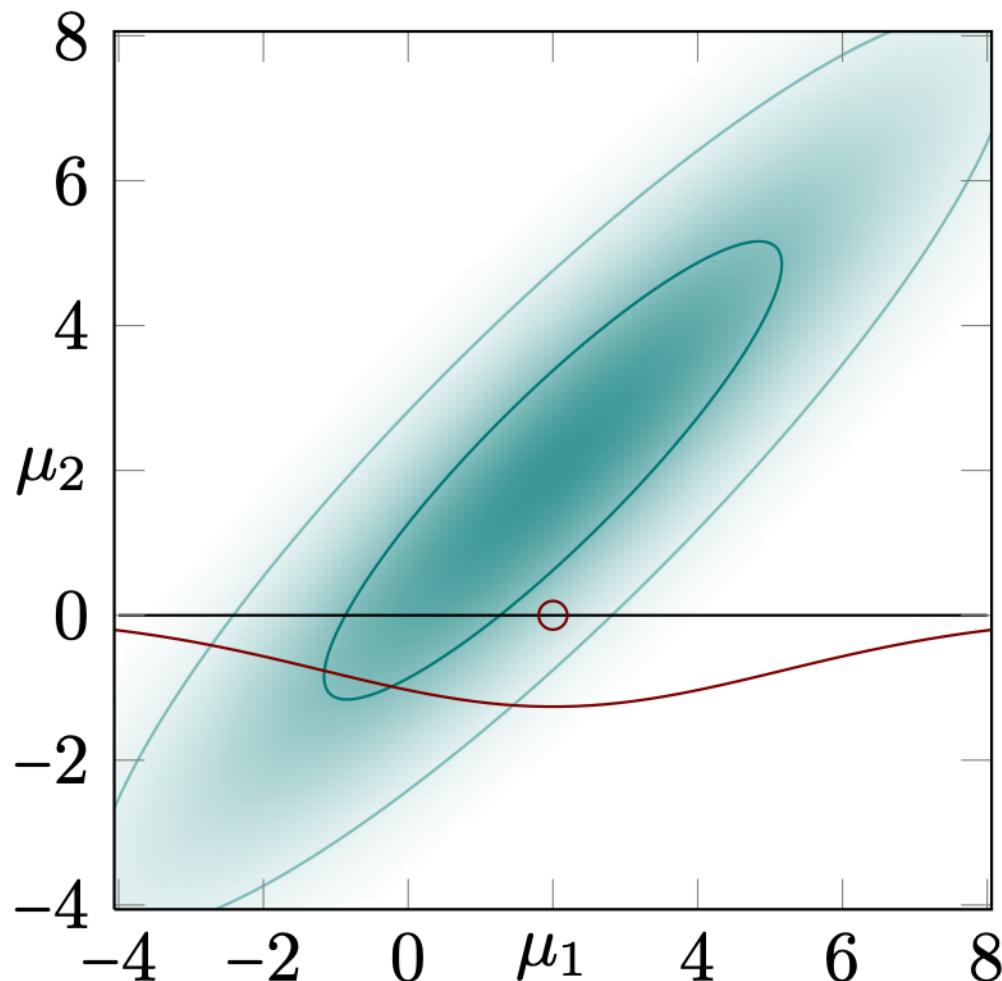
**Gaussians** are for **inference** what **linear** maps are for **algebra**.

# Closure under Marginalization

projections of Gaussians are Gaussian

- ▶ projection with  $A = \begin{pmatrix} 1 & 0 \end{pmatrix}$

$$\int \mathcal{N}\left[\begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$



- ▶ this is the **sum rule**

$$\int p(x, y) dy = \int p(y | x)p(x) dy = p(x)$$

- ▶ so every finite-dim Gaussian is a marginal of **infinitely many more**

# Gaussian process covariance functions (kernels)

$p(f)$  is a **Gaussian process** if for any finite subset  $\{x_1, \dots, x_n\} \subset \mathcal{X}$ , the marginal distribution over that finite subset  $p(f)$  has a multivariate Gaussian distribution.

Gaussian processes (GPs) are parameterized by a **mean function**,  $\mu(x)$ , and a **covariance function**, or **kernel**,  $K(x, x')$ .

$$p(f(x), f(x')) = \mathcal{N}(\mu, \Sigma)$$

where

$$\mu = \begin{bmatrix} \mu(x) \\ \mu(x') \end{bmatrix} \quad \Sigma = \begin{bmatrix} K(x, x) & K(x, x') \\ K(x', x) & K(x', x') \end{bmatrix}$$

and similarly for  $p(f(x_1), \dots, f(x_n))$  where now  $\mu$  is an  $n \times 1$  vector and  $\Sigma$  is an  $n \times n$  matrix.

# Gaussian process covariance functions

Gaussian processes (GPs) are parameterized by a mean function,  $\mu(x)$ , and a covariance function,  $K(x, x')$ .

An example covariance function:

$$K(x_i, x_j) = v_0 \exp \left\{ - \left( \frac{|x_i - x_j|}{r} \right)^\alpha \right\} + v_1 + v_2 \delta_{ij}$$

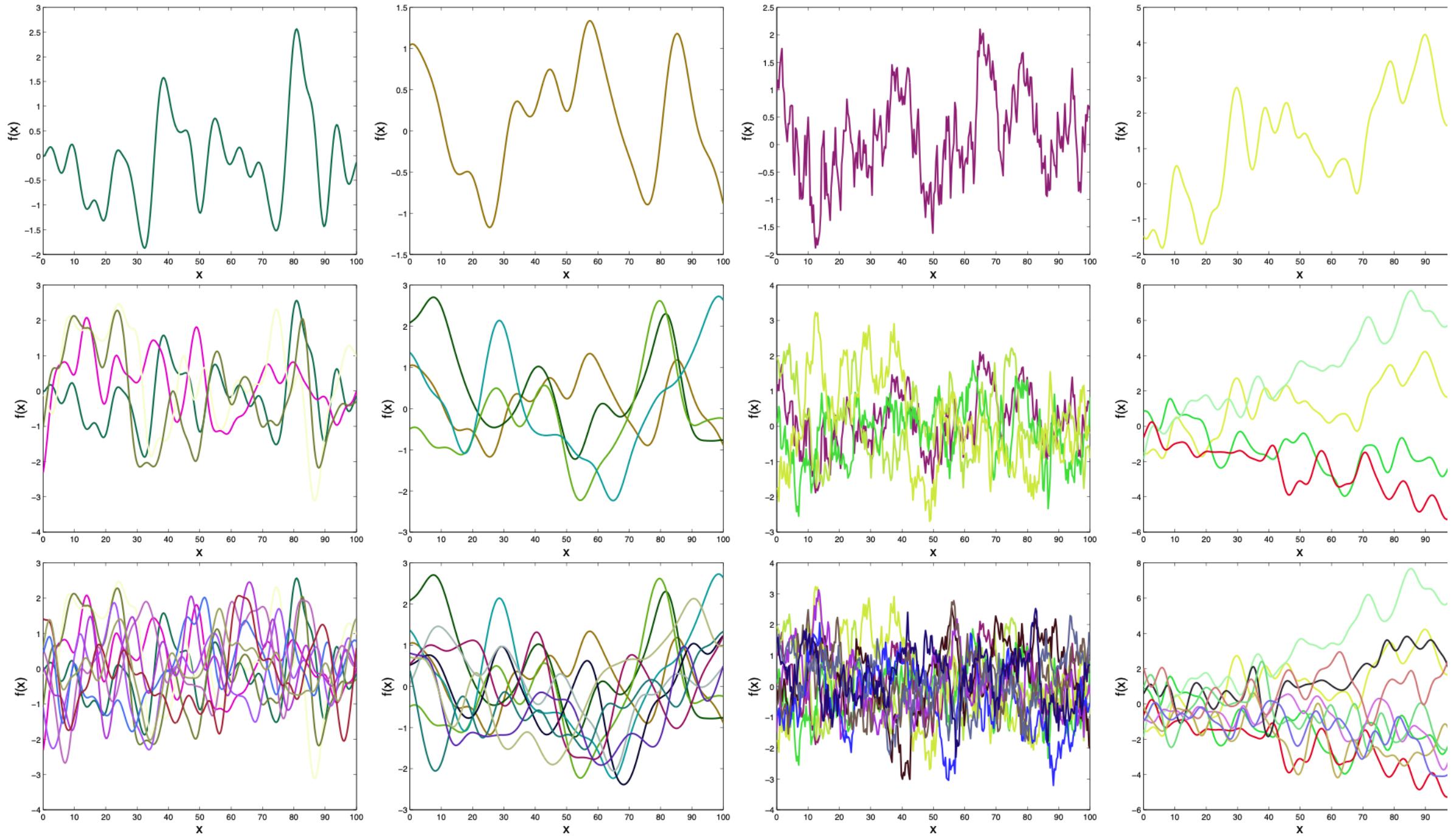
with parameters  $(v_0, v_1, v_2, r, \alpha)$

These kernel parameters are **interpretable** and can be learned from data:

$v_0$	signal variance
$v_1$	variance of bias
$v_2$	noise variance
$r$	lengthscale
$\alpha$	roughness

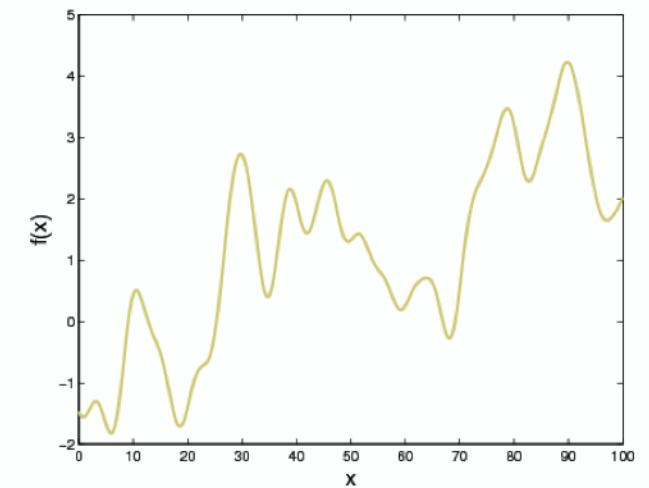
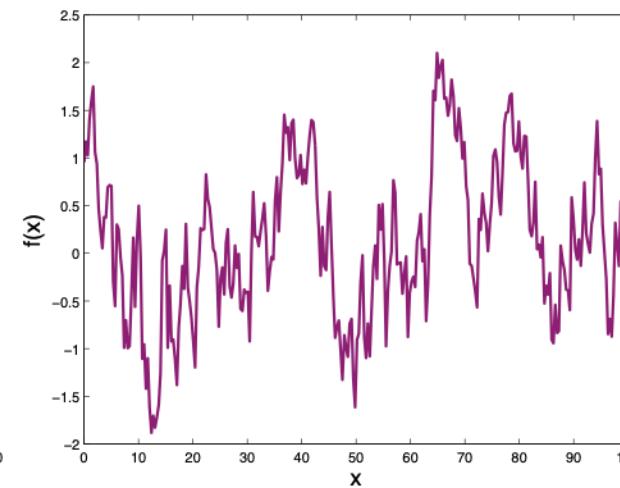
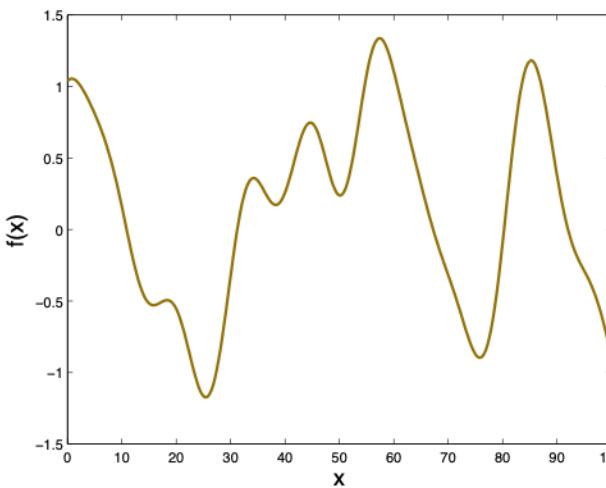
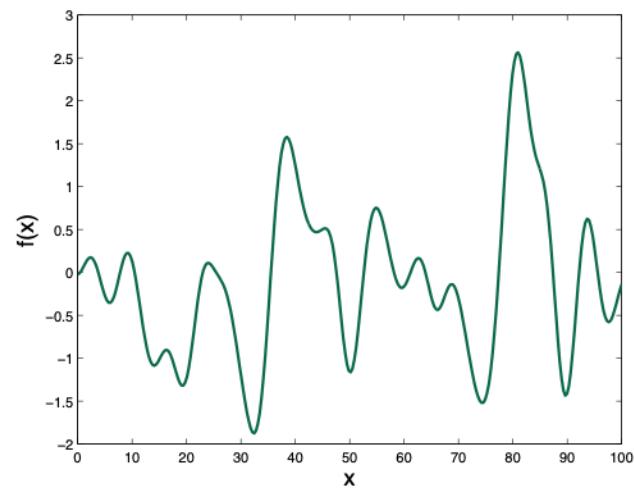
Once the mean and covariance functions are defined, everything else about GPs follows from the basic rules of probability applied to multivariate Gaussians.

# Samples from GPs with different $K(x, x')$

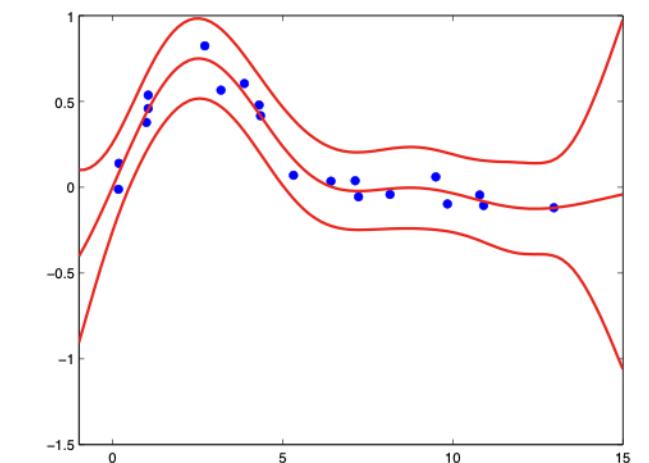
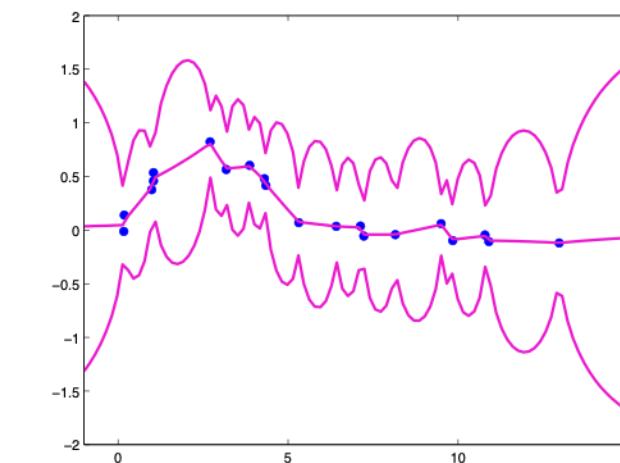
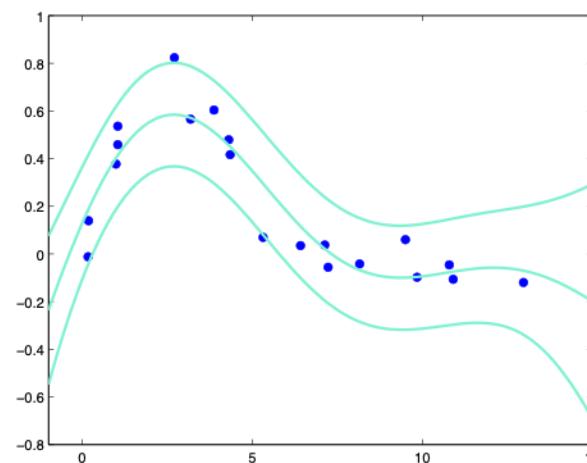
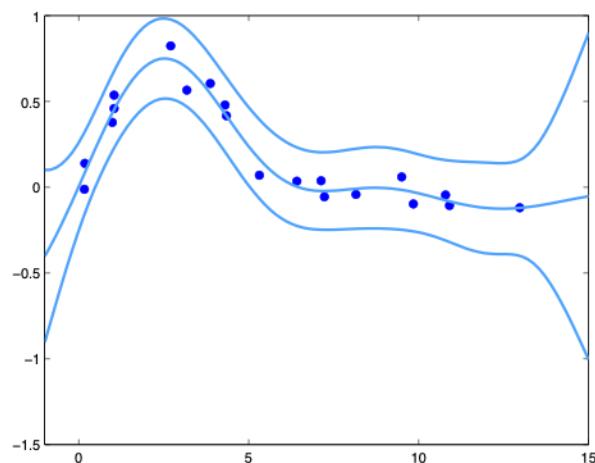


# Prediction using GPs with different $K(x, x')$

A sample from the prior for each covariance function:



Corresponding predictions, mean with two standard deviations:



# Using Gaussian processes for nonlinear regression

Imagine observing a data set  $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{i=1}^n\} = (\mathbf{X}, \mathbf{y})$ .

Model:

$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

$$f \sim \text{GP}(\cdot | 0, K)$$

$$\epsilon_i \sim \mathcal{N}(\cdot | 0, \sigma^2)$$

Prior on  $f$  is a GP, likelihood is Gaussian, therefore posterior on  $f$  is also a GP.

We can use this to make **predictions**

$$p(y_* | \mathbf{x}_*, \mathcal{D}) = \int p(y_* | \mathbf{x}_*, f, \mathcal{D}) p(f | \mathcal{D}) df$$

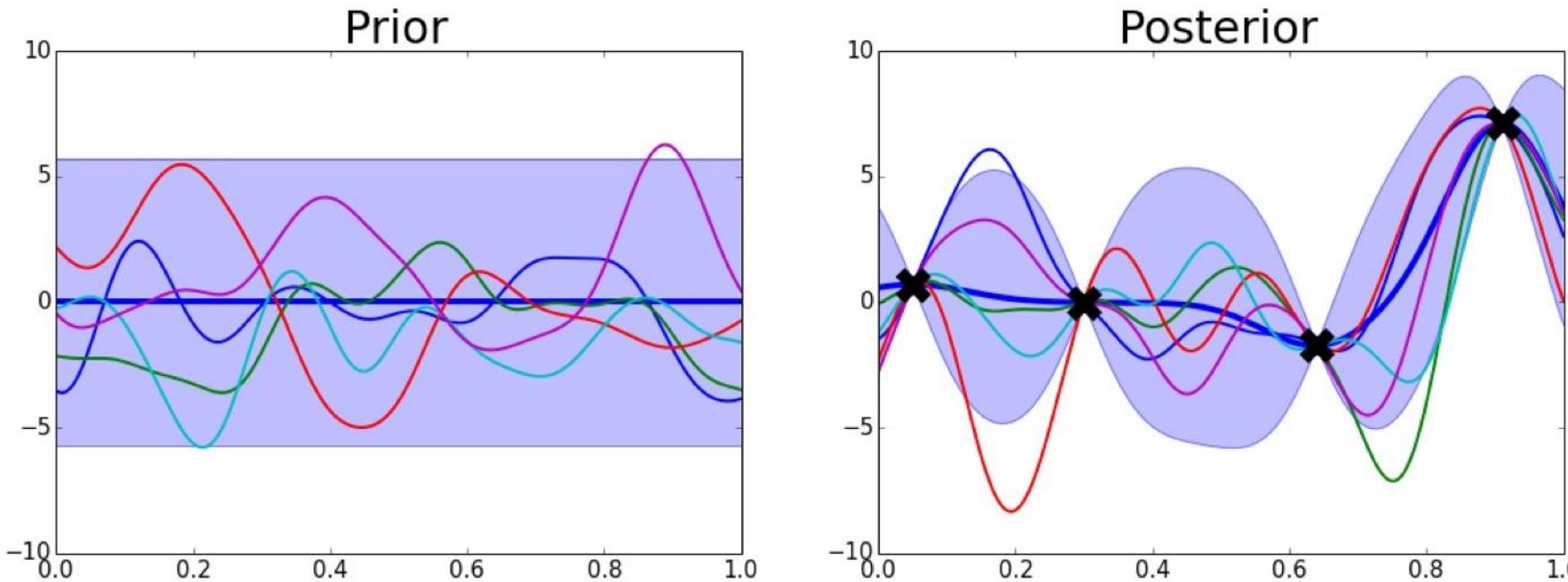
We can also compute the **marginal likelihood** (evidence) and use this to compare or tune covariance functions

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | f, \mathbf{X}) p(f) df$$

# Data-driven modeling with Gaussian processes

$$y = f(\mathbf{x}) + \epsilon$$

$$f \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}))$$



Training via maximizing the marginal likelihood

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\mathbf{K} + \sigma_\epsilon^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log 2\pi$$

Prediction via conditioning on available data

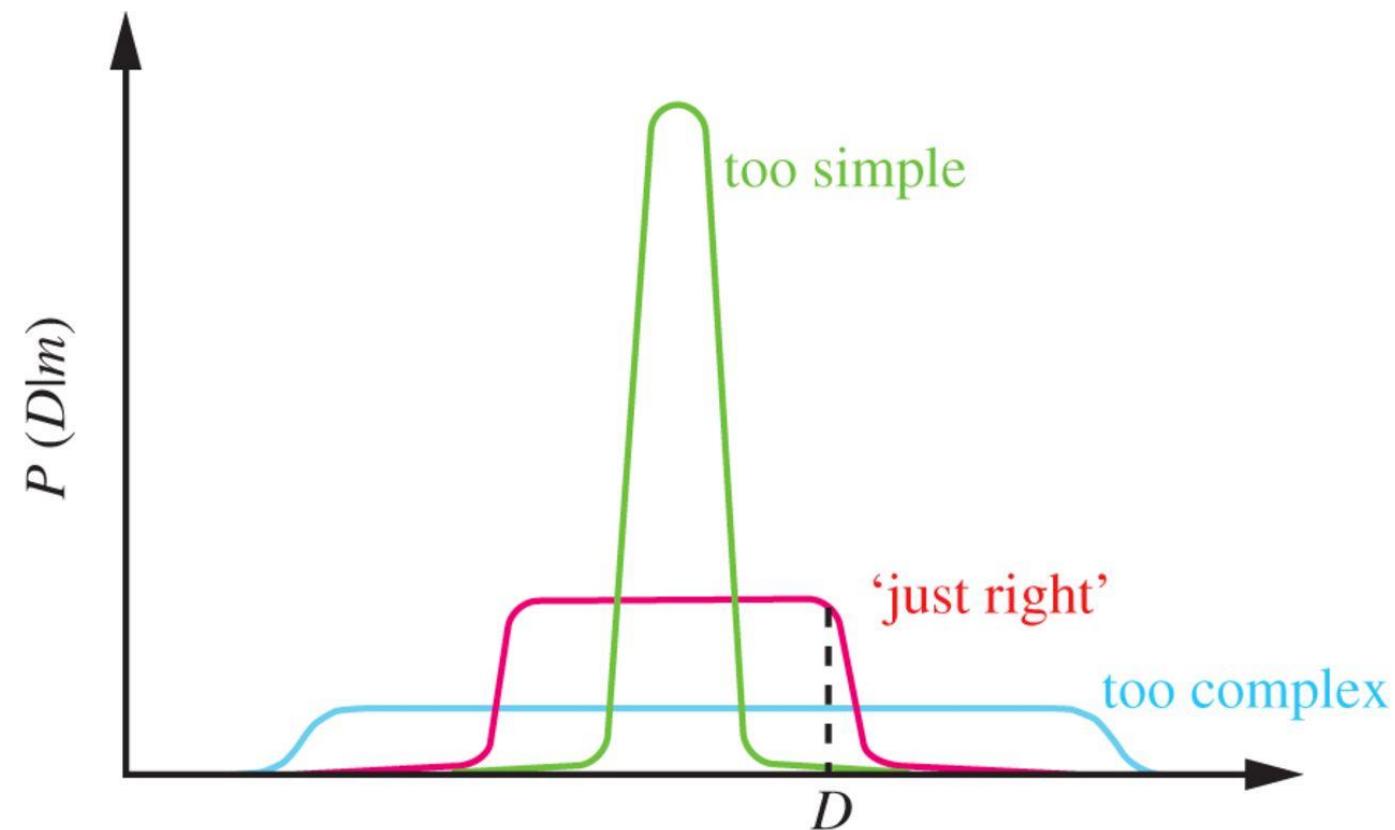
$$p(f_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$\mu_*(\mathbf{x}_*) = \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y},$$

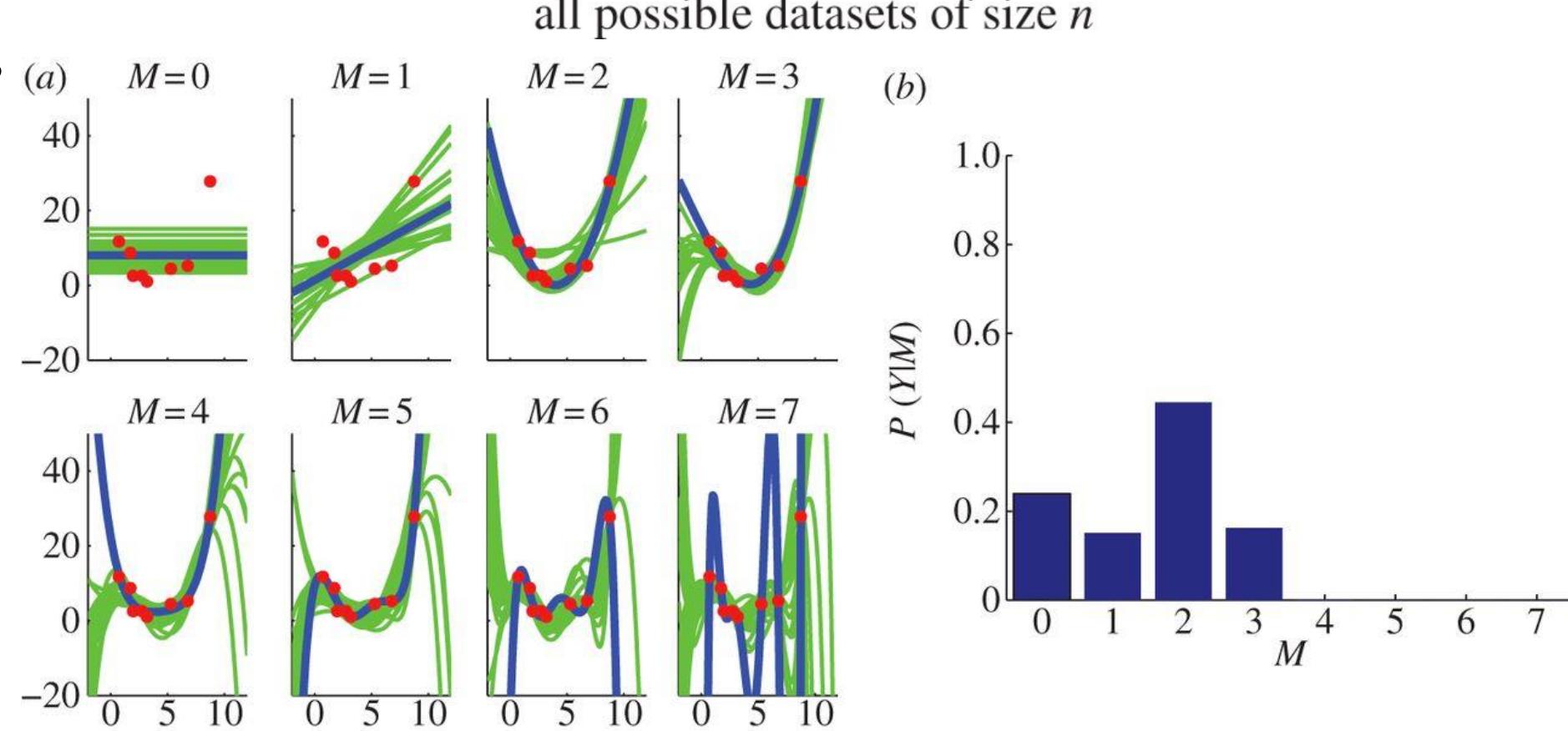
$$\sigma_*^2(\mathbf{x}_*) = \mathbf{k}_{**} - \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_{N*},$$

# Occam's razor

William of Ockham (~1285-1347 A.D)



**"plurality should not be posited without necessity."**



# Summary

## 1. Prior

$$\begin{aligned}y &= f(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2 I), x \in \mathbb{R}^d, y \in \mathbb{R} \\f(x) &\sim GP(0, \kappa(x, x'; \theta)) \\ \begin{bmatrix} f(x) \\ f(x') \end{bmatrix} &\sim GP \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \kappa(x, x) & \kappa(x, x') \\ \kappa(x', x) & \kappa(x', x') \end{bmatrix} \right)\end{aligned}$$

For example, consider the radial basis function:

$$\begin{aligned}\kappa(x, x'; \theta) &= \sigma_f^2 \exp \left( -\frac{1}{2} \sum_{i=1}^d \frac{(x_i - x'_i)^2}{\theta_i^2} \right) \\ \Theta := \{\theta, \sigma_n^2\} &= \{\sigma_f^2, \theta_1, \dots, \theta_d, \sigma_n^2\}\end{aligned}$$

## 2. Training

Given data  $\{X, y\}$ ,  $X \in \mathbb{R}^{N \times d}, y \in \mathbb{R}^{N \times 1}$

$$\begin{aligned}p(y|x) &= \mathcal{N}(0, \underbrace{\kappa(x, x; \theta) + \sigma_n^2 I}_K) \\ \implies \log(p(y|x)) &= \frac{1}{2} y^T K^{-1} y + \frac{1}{2} \log |K| + \frac{n}{2} \log(2\pi)\end{aligned}$$

Determine parameters to maximize log-likelihood.  $K$  is  $n \times n$ , full, symmetric positive definite. Can be solved in  $\mathcal{O}(N^3)$ . Gradients  $\nabla_\theta \log p(y|x)$  can be obtained other analytically or through autodiff.

## 3. Prediction / Posterior

$$\begin{aligned}\begin{bmatrix} f(x^*) \\ y \end{bmatrix} &\sim GP \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \kappa(x^*, x^*) & \kappa(x^*, x) \\ \kappa(x, x^*) & \underbrace{\kappa(x, x)}_K \end{bmatrix} \right) \\ \implies p(f(x^*|x, y) &= \mathcal{N} \left( \underbrace{\kappa(x^*, x) K^{-1} y}_{\mu(x^*)}, \underbrace{\kappa(x^*, x^*) - \kappa(x^*, x) K^{-1} \kappa(x, x^*)}_{\Sigma(x^*, x^*)} \right)\end{aligned}$$