

Prop:  $\underline{X} \sim N_m(0, \Sigma)$   $i, j \in \{1, \dots, m\}$

$$B = \{1, \dots, m\} \setminus \{i, j\}$$

$$A = \{i, j\}$$

TFAE

$$\text{Var}(X_A | X_B) =$$

$$= \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

$$(i) X_i \perp\!\!\!\perp X_j | X_B$$

$$(ii) \Sigma_{ij} = \underbrace{\Sigma_{iB}}_{\substack{\uparrow \\ \mathbb{R}^{m-2}}} \Sigma_{BB}^{-1} \Sigma_{Bj}$$

$$(iii) \left( \text{Cov}(X_i, X_j | X_B) = 0 \right)$$

$$(iii) (\Sigma^{-1})_{ij} = 0$$

$$(iv) \text{the } j^{\text{th}} \text{ entry of } \Sigma_{i, B \cup \{j\}}^{-1} \Sigma_{B \cup \{j\}, B \cup \{j\}}$$

is zero,  $\substack{\uparrow \\ \mathbb{R}^{m-1}}$

DATA

$$\underline{X}^{(1)}, \dots, \underline{X}^{(n)} \sim N_m(\mu, \Sigma)$$

$\mu, \Sigma$  unknown

$$f(x) = \frac{1}{(2\pi)^{m/2}} (\det \Sigma)^{-1/2} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\log f(x) = -\frac{m}{2} \log(2\pi) - \frac{1}{2} \log \det \Sigma - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$$

log-likelihood

$$l(\mu, \Sigma) = \sum_{i=1}^n \log f(x^{(i)}) =$$

$$= -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log \det \Sigma$$

$$- \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)$$

→ optimize wot  $\mu \in \mathbb{R}^m$

$$(x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu) =$$

$$x^{(i)T} \Sigma^{-1} x^{(i)} \textcircled{1}$$

equal

$$- x^{(i)T} \Sigma^{-1} \mu \textcircled{2}$$

$$- \mu^T \Sigma^{-1} x^{(i)} \textcircled{3}$$

$$+ \mu^T \Sigma^{-1} \mu$$

$$\textcircled{1} \quad (\Sigma^{-1} x^{(i)})^T \mu$$

$$\textcircled{2} \quad \mu^T (\Sigma^{-1} x^{(i)})$$

Note  $w \in \mathbb{R}^m$   $y \in \mathbb{R}^m$

$$\nabla_y (w^T y) = \nabla_y (y^T w) = w$$

$$\frac{\partial}{\partial y_k} (w^T y) = \frac{\partial}{\partial y_k} (w_1 y_1 + \dots + w_m y_m) = w_k$$

② A symmetric

$$\nabla_y (y^T A y) = 2 \cdot A \cdot y$$

$$\frac{\partial}{\partial y_k} (y^T A y) \stackrel{?}{=} (2A y)_k$$

$$\stackrel{?}{=} \frac{\partial}{\partial y_k} \left( \sum_{i,j=1}^m A_{ij} y_i y_j \right)$$

$$\stackrel{?}{=} \frac{\partial}{\partial y_k} \left( \sum_{i=1}^m A_{ii} y_i^2 + 2 \cdot \sum_{i < j} A_{ij} y_i y_j \right)$$

$$\begin{aligned}
&= 2A_{kk}y_k + 2 \underbrace{\sum_{i \neq k} A_{ik}y_i}_{\hookrightarrow A \text{ symmetric}} \\
&= 2 \sum_{i=1}^m A_{ik}y_i = 2 \sum_{i=1}^m A_{ki}y_i \\
&= (2Ay)_k
\end{aligned}$$


---

$$\begin{aligned}
&\nabla \left( (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu) \right) = \\
&= \nabla \textcircled{1} - \nabla \textcircled{2} - \nabla \textcircled{3} + \nabla \textcircled{4} \\
&= 0 - \Sigma^{-1} x^{(i)} - \Sigma^{-1} x^{(i)} \\
&\quad + 2 \cdot \Sigma^{-1} \mu \\
&= 2 \Sigma^{-1} (\mu - x^{(i)})
\end{aligned}$$

$$\nabla_{\mu}(\ell(\mu, \Sigma)) = -\frac{1}{2} \sum_{i=1}^n \left( 2 \Sigma^{-1} (\mu - x^{(i)}) \right)$$

$$= \sum_{i=1}^n \Sigma^{-1} (x^{(i)} - \mu)$$

$$= \Sigma^{-1} \sum_{i=1}^n (x^{(i)} - \mu) \stackrel{?}{=} 0$$

irrespective of what  $\Sigma$  is  
the solution is the same  
as the solution to

$$\sum_{i=1}^n (x^{(i)} - \mu) = 0$$

equiv.  $\frac{1}{n} \sum_{i=1}^n x^{(i)} - \mu = 0$

MLE  $\hat{\mu} = \bar{X}_n$

Recall:  $\bar{X}_n \sim N(\mu, \frac{1}{n})$

plug  $\mu = \hat{\mu}$

$$l(\bar{X}_n, \Sigma) = -\frac{nm}{2} \log(2\pi)$$

$$- \frac{n}{2} \log \det \Sigma$$

$$- \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \bar{X}_n)^T \Sigma^{-1} (x^{(i)} - \bar{X}_n)$$

$$\left\{ \begin{array}{l} A \in \mathbb{R}^{n \times m} \quad B \in \mathbb{R}^{m \times n} \end{array} \right.$$

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\left\{ \begin{array}{l} \underline{x}, \underline{y} \in \mathbb{R}^m \end{array} \right.$$

$$\underline{x}^T \underline{y} = \text{tr}(\underline{x}^T \underline{y}) = \text{tr}(\underline{y} \underline{x}^T)$$

$$l(\bar{X}_n, \Sigma) = -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log \det \Sigma$$
$$- \frac{1}{2} \sum_{i=1}^n \text{tr}(\Sigma^{-1} (x^{(i)} - \bar{X}_n) (x^{(i)} - \bar{X}_n)^T)$$

$$= -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log \det \Sigma$$

$$- \frac{n}{2} \operatorname{tr} \left( \Sigma^{-1} \left( \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{x}_n) (x^{(i)} - \bar{x}_n)^T \right) \right)$$

$S_n$

$$= -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log \det \Sigma$$

$$- \frac{n}{2} \operatorname{tr} (\Sigma^{-1} S_n)$$

param.  $K = \Sigma^{-1}$

$$l(\bar{x}; K) = -\frac{nm}{2} \log(2\pi) + \frac{n}{2} \log \det K$$

$\uparrow$  inv. cov

$$- \frac{n}{2} \operatorname{tr} (K S_n)$$

$\nabla_K$  matrix with entries  $\frac{\partial l}{\partial K_{ij}}$

$$\nabla_K l(\bar{x}; K) = \frac{n}{2} K^{-1} - \frac{n}{2} S_n$$

$$\text{MLE } \hat{K} : \hat{K}^{-1} = S_n$$

- is  $S_n$  not invertible, no solution
- if it is equiv.  $\hat{\Sigma} = S_n$

---

Prop :  $\underline{x}^{(1)}, \dots, \underline{x}^{(n)} \sim N(\mu, \Sigma)$

$$\bar{X}_n \perp S_n \quad \underline{X} = \begin{pmatrix} \underline{x}^{(1)} \\ \vdots \\ \underline{x}^{(n)} \end{pmatrix}$$

Proof : Recall  $S_n = \frac{1}{n} X^T H X$

$$H = I_n - \frac{1}{n} \underline{1}_n \underline{1}_n^T \quad \text{centering matrix}$$

$$S_n = \frac{1}{n} (HX)^T (HX)$$

equiv. show  $\underline{\bar{X}}_n \perp \underline{HX}$

$$k^{\text{th}} \text{ row of } HX = x^{(k)} - \bar{X}_n$$



$$\begin{aligned}
& \text{Cov}(\bar{X}_n, X^{(k)} - \bar{X}_n) \\
&= \text{Cov}\left(\bar{X}_n, X^{(k)}\right) - \underbrace{\text{Var}(\bar{X}_n)}_{\frac{1}{n} \cdot \Sigma} \\
&= \frac{1}{n} \cdot \underbrace{\text{Cov}(X^{(k)}, X^{(k)})}_{\Sigma} - \frac{1}{n} \Sigma = 0_{m \times m}
\end{aligned}$$

so  $HX \perp \bar{X}$

so  $S_n \perp \bar{X}$



STA 437/2005:  
Methods for Multivariate Data  
Week 4: Gaussian Processes

Piotr Zwiernik

University of Toronto

# Table of contents

1. Introduction to Gaussian Processes (GPs)
2. GPs for Spatial Data
3. Nonparametric Regression with GPs

# Introduction to GPs

# Marginal distribution of MVN

$$X = (X_1, \dots, X_m)$$

Consider the following reformulation of the earlier result:

Suppose  $X \sim N_m(\mu, \Sigma)$ . Let  $T := \{1, \dots, m\}$  and define  $\mu_i = \mathbb{E}X_i$

- ▶  $m : T \rightarrow \mathbb{R}$  such that  $m(i) := \mu_i$  (mean function)
- ▶  $k : T \times T \rightarrow \mathbb{R}$  such that  $k(i, j) := \Sigma_{ij}$  (kernel function)

Then for every  $A = \{t_1, \dots, t_n\} \subseteq T$ , the vector  $X_A = (X_{t_1}, \dots, X_{t_n})$  is Gaussian with

- ▶ The mean  $\mu_A$  whose  $i$ -th entry is  $m(t_i)$ .
- ▶ The covariance matrix  $\Sigma_{AA}$  whose  $(i, j)$ -th entry is  $k(t_i, t_j)$ .

The set  $T$  indexes all random variables in the system.

For every  $A = \{t_1, \dots, t_n\} \subseteq T$ ,  $(X_{t_1}, \dots, X_{t_n})$  is Gaussian.

# Gaussian Processes - an immediate generalization

A **Gaussian Process (GP)** is a generalization of the multivariate normal distribution to a collection of random variables indexed by an **arbitrary** set  $T$ .

## Definition

A Gaussian Process is a collection of random variables  $\{X_t\}_{t \in T}$  such that for any finite set of points  $\{t_1, \dots, t_n\} \subset T$ , the corresponding vector  $(X_{t_1}, \dots, X_{t_n})$  follows a multivariate normal distribution.

In what follows we assume  $T \subseteq \mathbb{R}^d$  with the Euclidean distance metric.

Often, the correlation between two variables  $X_s$  and  $X_t$  will depend on the distance  $\|t - s\|$ .

# The mean and the kernel functions

A Gaussian Process is characterized by:

- ▶ A **mean function**  $m : T \rightarrow \mathbb{R}$ :  $m(t) = \mathbb{E}[X_t]$
- ▶ A **kernel function**  $k : T \times T \rightarrow \mathbb{R}$ :  $k(t, t') = \text{Cov}(X_t, X_{t'})$

Note that  $m$  is pretty much arbitrary (often set to be zero) but  $k$  is highly constrained:

Positive semi-definiteness:

For any finite set  $\{t_1, \dots, t_n\} \subset T$ , the covariance matrix  $\Sigma$  with entries  $\Sigma_{ij} = k(t_i, t_j)$  is positive semi-definite.

We can use feature maps  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^p$  to define kernels:

$$k(s, t) = \psi(s)^\top \psi(t).$$

Feature maps define kernels but not all kernels are like that (this can be generalized to “infinite dimensional” feature maps).

# Common Kernels in GPs

## ► Squared Exponential (RBF) Kernel:

$$k_E(t, t') = \sigma^2 \exp\left(-\frac{\|t - t'\|^2}{2\ell^2}\right).$$

- Controls smoothness of the functions sampled from the GP.
- Length scale  $\ell$ : Correlation distance.
- Signal variance  $\sigma^2$ : Scale of the output.

## ► Matérn Kernel:

$$k_M(t, t') = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \frac{\|t - t'\|}{\ell}\right)^\nu K_\nu\left(\sqrt{2\nu} \frac{\|t - t'\|}{\ell}\right).$$

- $\nu$ : Smoothness parameter.
- More flexible than the RBF kernel for modeling rough functions.



# Constructing kernels from kernels

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , the following kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad \text{for } c > 0,$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top A \mathbf{x}' \quad (A \text{ PSD})$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$$

where  $q$  polynomial with  $\geq 0$  coefficients.

# Modelling with Gaussian processes

Working with Gaussian Processes we fix a kernel function.

Data: Suppose we observed  $(X_{t_1}, \dots, X_{t_n})$  for some  $t_1, \dots, t_n \in T$ .

If the kernel function comes with some hyperparameters  $\alpha$ , we can learn them maximizing the log-likelihood.

- ▶ By definition,  $(X_{t_1}, \dots, X_{t_n})$  is MVN with covariance that depends on  $\alpha$ .
- ▶ This may be a complicated optimization procedure.

Suppose we want to predict the value of the process at some point  $t_{n+1}$

- ▶ By definition  $(X_{t_1}, \dots, X_{t_n}, X_{t_{n+1}})$  is jointly Gaussian so simply compute the conditional distribution:  $X_{t_{n+1}} | X_{t_1}, \dots, X_{t_n}$ .
- ▶ This gives both the point prediction (the conditional mean) and uncertainty quantification (conditional variance).

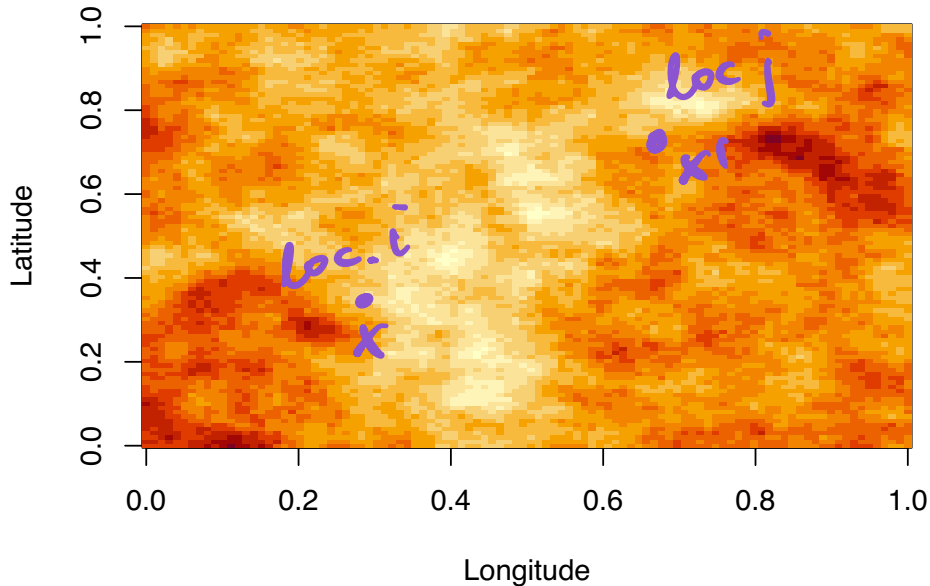
# GPs for Spatial Data

# Example: Modeling Spatial Data with GPs

$$T = [0, 1]^2$$

GPs are widely used in spatial statistics, e.g. temperature across a grid of locations.

$$\underline{x}, \underline{x}' \in T$$



- Grid of  $100^2$  points.
- Fix the exponential kernel  $\exp\{-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|\}$
- Compute the  $100^2 \times 100^2$  covariance matrix
- Get 1 sample from the corresponding distr.

Handling a 10000-dimensional Gaussian comes with its own computational challenges.

# Spatial GP: Prediction

We explained how to make a prediction for  $X_{t_{n+1}}$ . This easily generalizes.

Suppose we observed the mean zero GP over some locations  $\mathbf{x}_{\text{train}}$ .

Our goal is to make predictions over some other points  $\mathbf{x}_{\text{test}}$

1. Combine training and test locations.
2. Compute the covariance matrix using the kernel function.
3. Use Gaussian conditioning formulas:

$$\begin{aligned}\mathbb{E}[\mathbf{x}_{\text{test}}|\mathbf{x}_{\text{train}}] &= \Sigma_{\text{test,train}}\Sigma_{\text{train,train}}^{-1}\mathbf{x}_{\text{train}}, \\ \text{COV}(\mathbf{x}_{\text{test}}|\mathbf{x}_{\text{train}}) &= \Sigma_{\text{test,test}} - \Sigma_{\text{test,train}}\Sigma_{\text{train,train}}^{-1}\Sigma_{\text{test,train}}.\end{aligned}$$

# Nonparametric Regression with GPs

# Nonparametric Regression

GPs can be used for nonparametric regression:

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

Prior over  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ : GP defined by  $m(\mathbf{x})$  and  $k(\mathbf{x}, \mathbf{x}')$ .

- ▶ In this sense GP defines a distribution over (random) functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ .

We have  $(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- ▶  $\mu_i = m(\mathbf{x}_i)$
- ▶  $\Sigma_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$

Say  $d = 1$ . Given  $m(x)$  and  $k(x, x')$ , how would you plot random samples of the corresponding random functions on  $\mathbb{R}$ ?

# Nonparametric Regression

Note that  $\mathbf{y} = (y_1, \dots, y_n) = (f(\mathbf{x}_1) + \varepsilon_1, \dots, f(\mathbf{x}_n) + \varepsilon_n)$ .

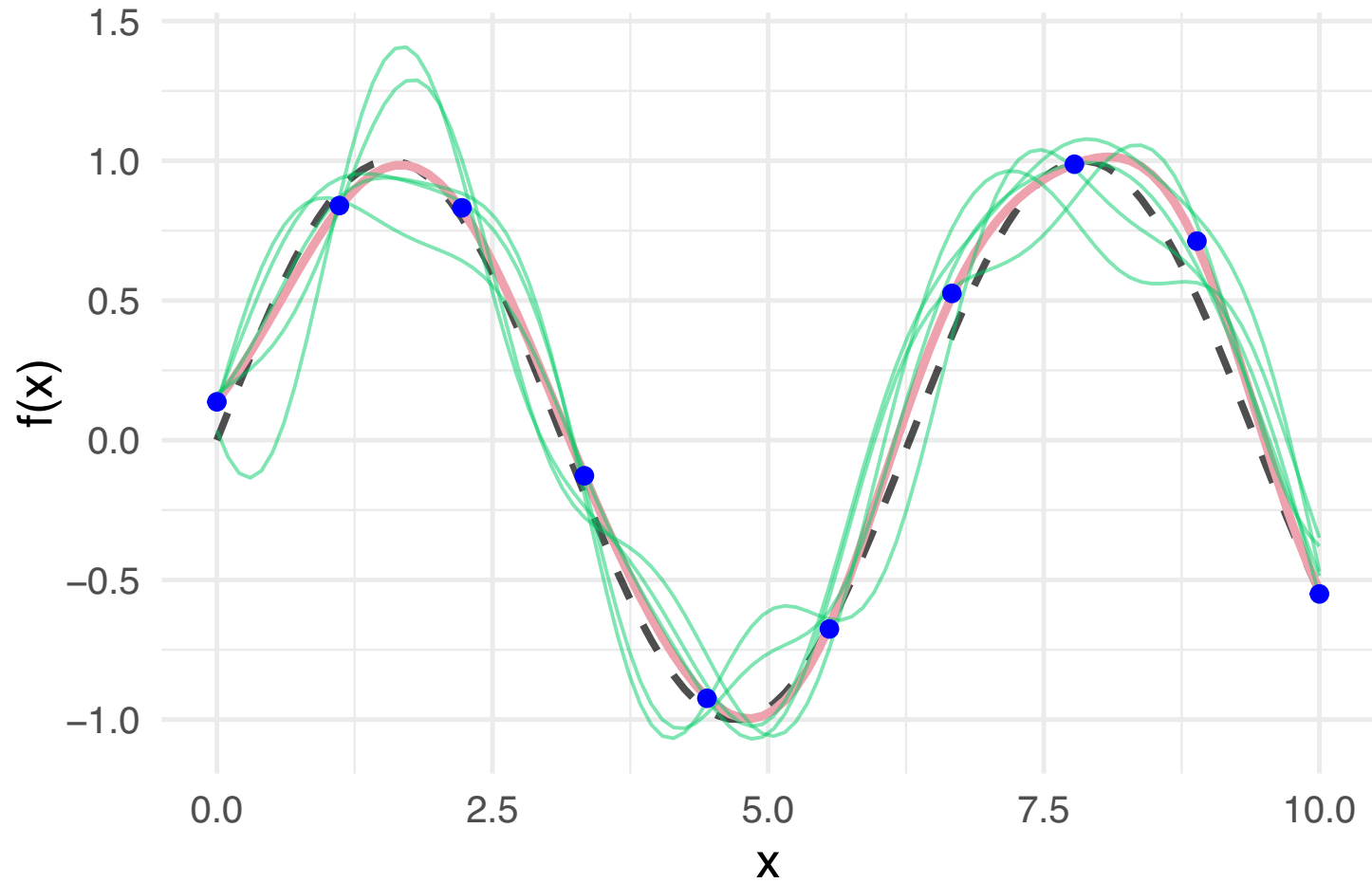
Consider the underlying Gaussian Process  $y(\mathbf{x})$ :

- The mean is  $m(\mathbf{x})$ .
  - ▶  $\mathbb{E}[y(\mathbf{x}_i)] = \mathbb{E}[f(\mathbf{x}_i) + \varepsilon_i] = m(\mathbf{x}_i)$ .
- The kernel is  $k(\mathbf{x}, \mathbf{x}') + \sigma^2 \mathbf{1}\{\mathbf{x} = \mathbf{x}'\}$ .
  - ▶  $\text{cov}[y(\mathbf{x}_i), y(\mathbf{x}_j)] = \text{cov}(f(\mathbf{x}_i) + \varepsilon_i, f(\mathbf{x}_j) + \varepsilon_j) = k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \mathbf{1}\{\mathbf{x}_i = \mathbf{x}_j\}$ .

Given data  $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$  we can now easily predict  $y$  at any other point  $\mathbf{x}$ .



## Gaussian Process Regression



# Summary

- ▶ Gaussian Processes are a versatile tool for regression and spatial modeling.
- ▶ Key components:
  - ▶ Mean function.
  - ▶ Kernel function.
- ▶ Takeaway: Conceptually it is not harder than MVNs and the same formulas apply.
- ▶ Computational issues can be significant.