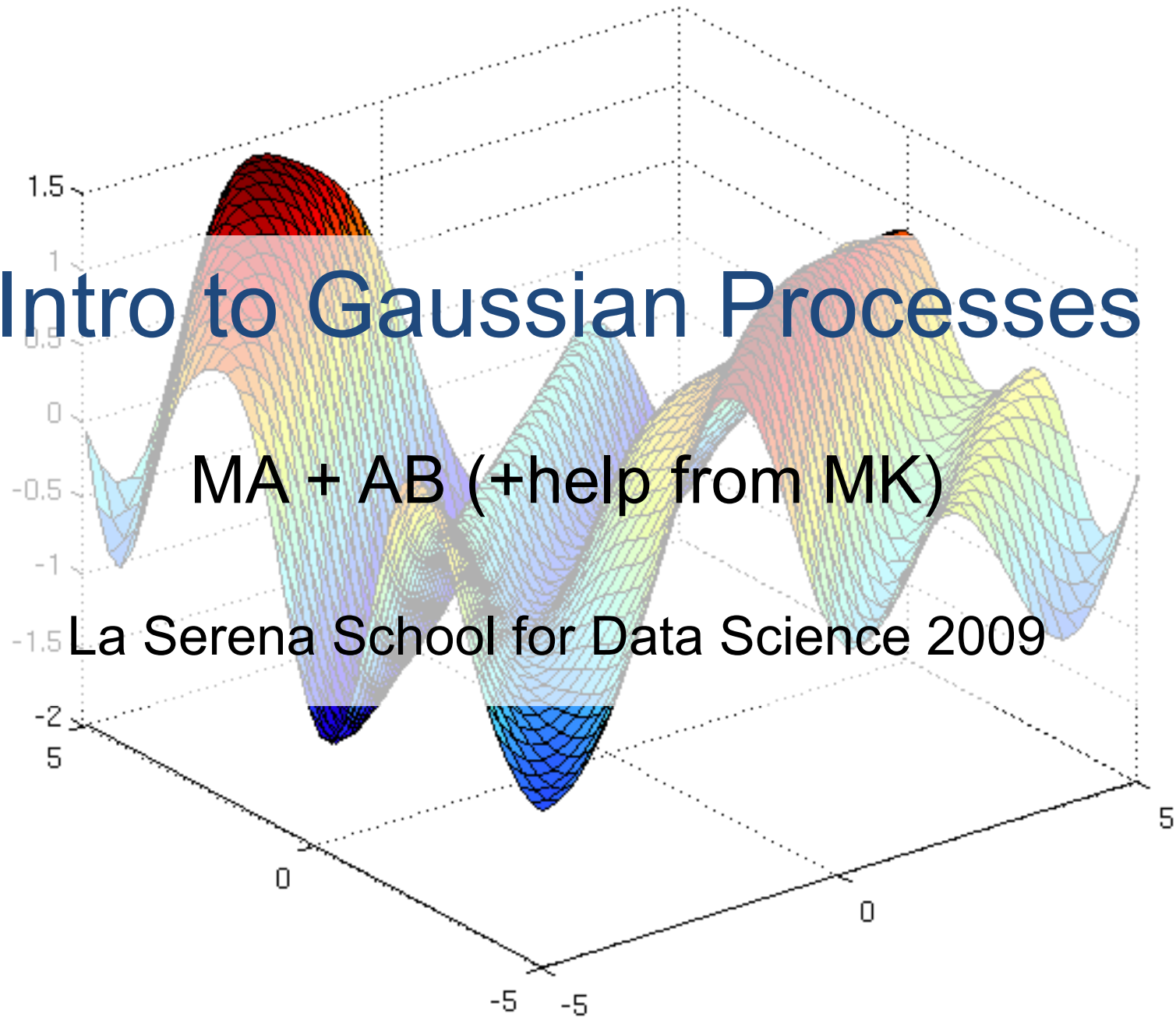


# Intro to Gaussian Processes

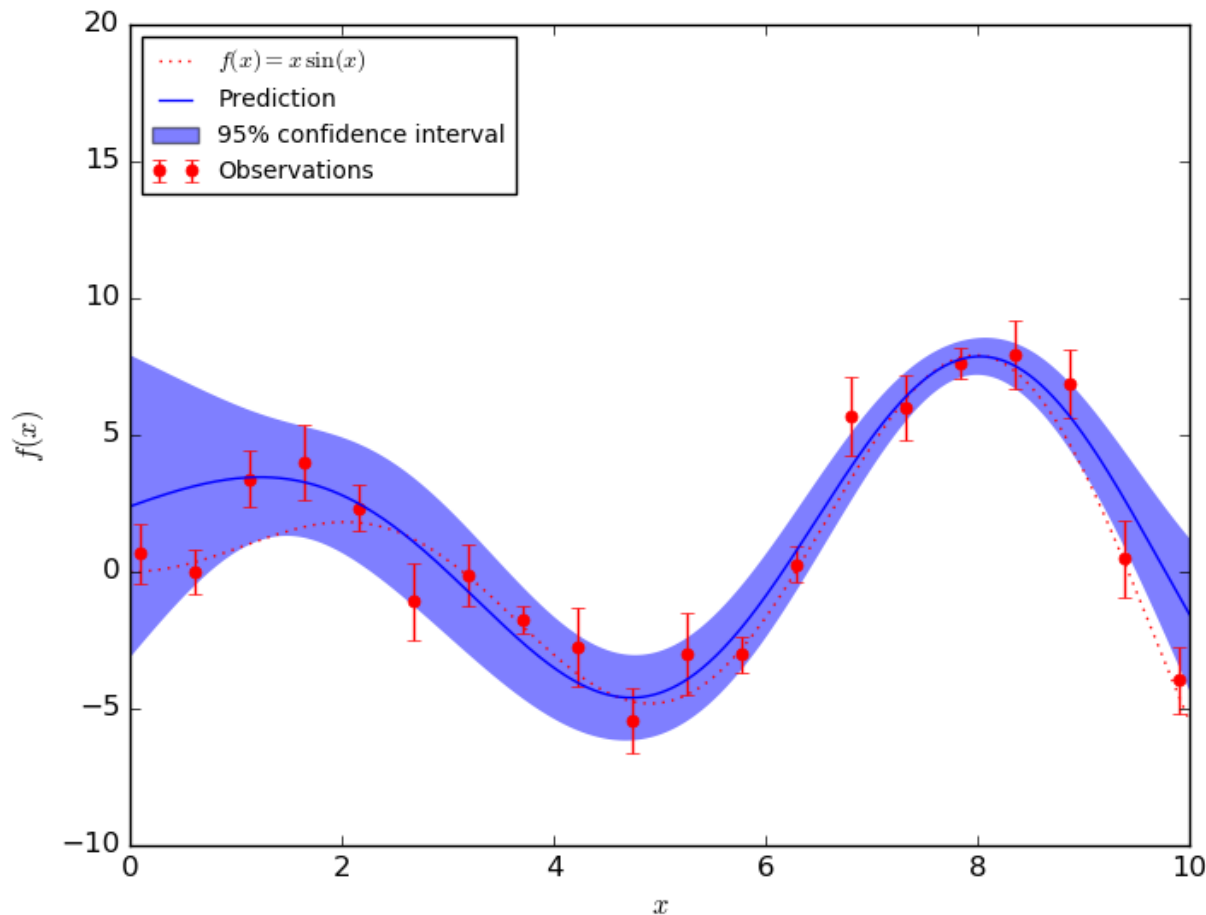
MA + AB (+help from MK)

La Serena School for Data Science 2009



# Motivation GP: non-parametric regression method

Stochastic process: generalization of pdf to “functions of pdfs” (if the process is “Gaussianly behaved (in time)” things “become easy”)



Bridging the gap  
between:

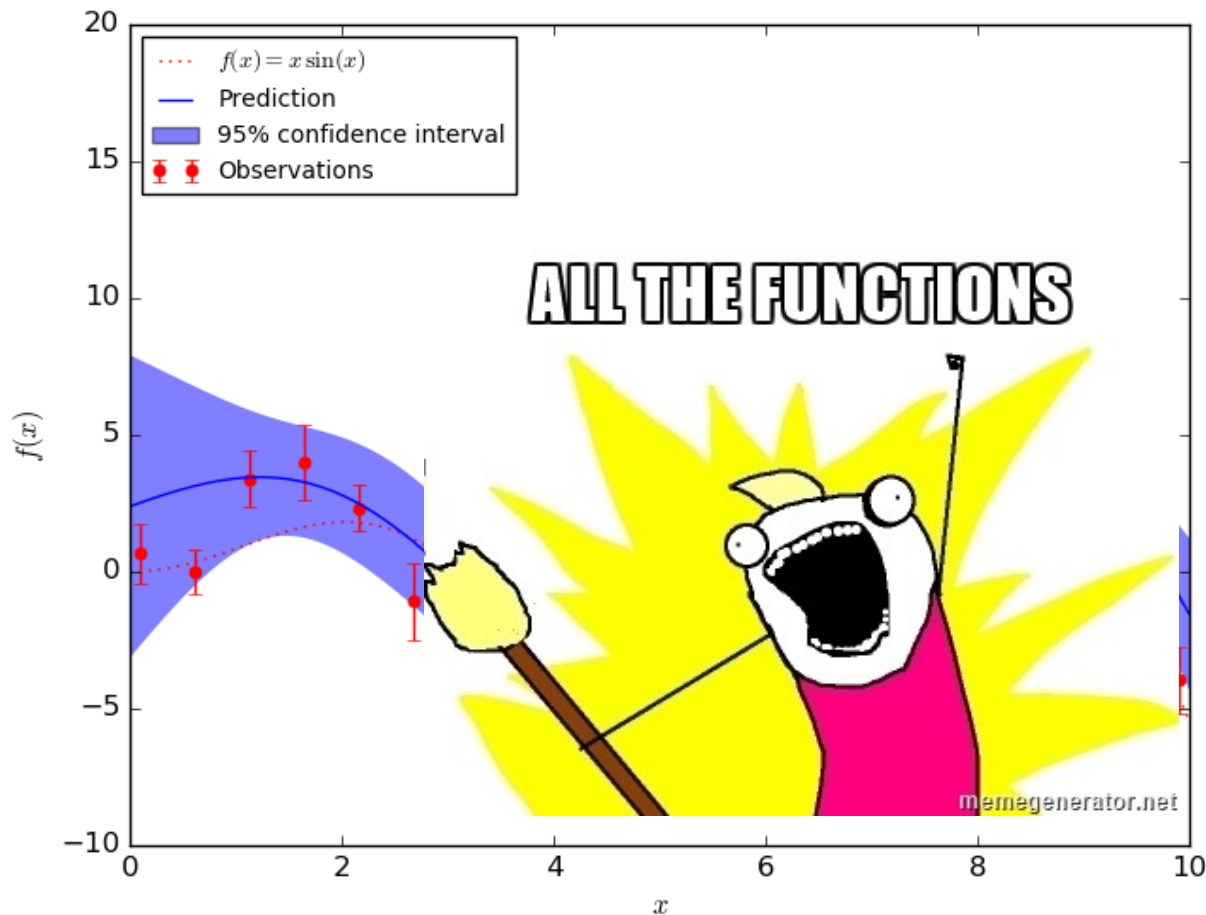
understanding data  
relationships

and

making predictions

# Motivation GP: non-parametric regression method

Stochastic process: generalization of pdf to “functions of pdfs” (if the process is “Gaussianly behaved (in time)” things “become easy”)



Bridging the gap  
between:

understanding data  
relationships

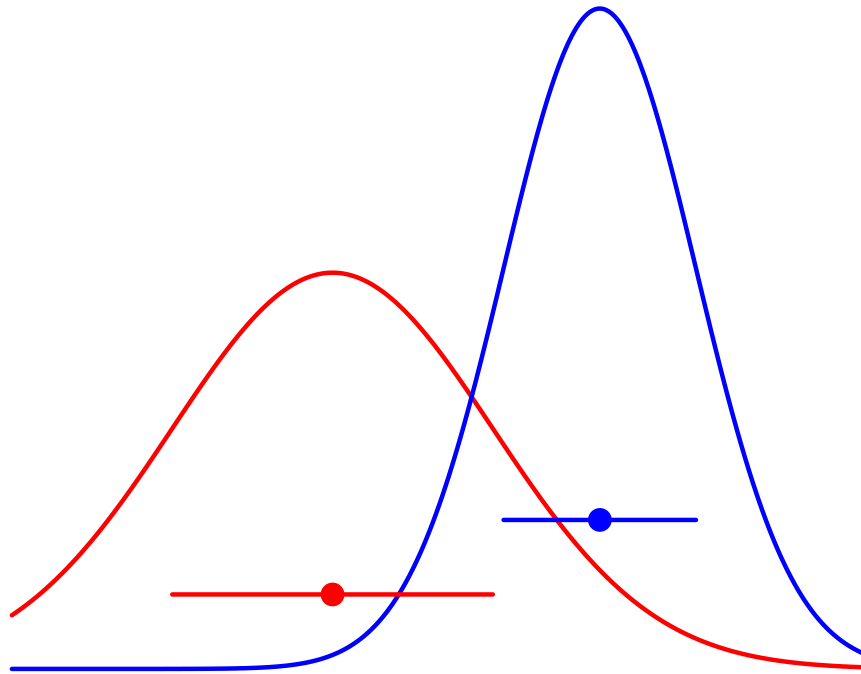
and

making predictions

# Motivation: the basics

---

A Gaussian distribution is described by a location  $\mu$  and “a shape”:  $\sigma$  for one dimension and for several, a covariance matrix  $\Sigma$

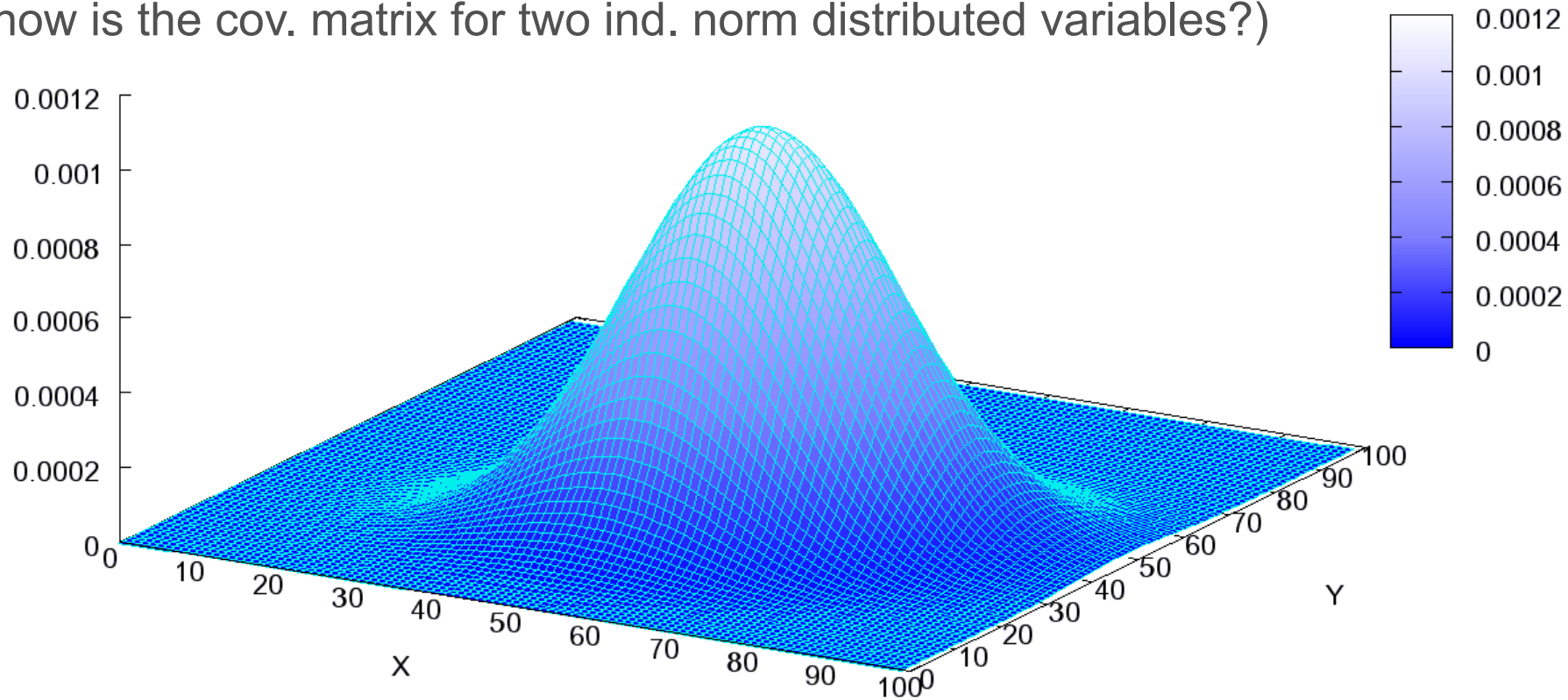


# Motivation: the basics

---

Multivariate Normal Distribution

Covariance matrix  $\rightarrow$  shape of the bell  
(how is the cov. matrix for two ind. norm distributed variables?)



# Motivation: the basics II

---

$$p(y|\theta) = \phi(y|\theta) = \phi(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

Single variable

$$p(\mathbf{Y}|\boldsymbol{\mu}, K) = \phi(\mathbf{Y}|\boldsymbol{\mu}, K) = \frac{\exp\left(-\frac{1}{2}(\mathbf{Y} - \boldsymbol{\mu})^T K^{-1}(\mathbf{Y} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^d \det K}}$$

Multi-variate

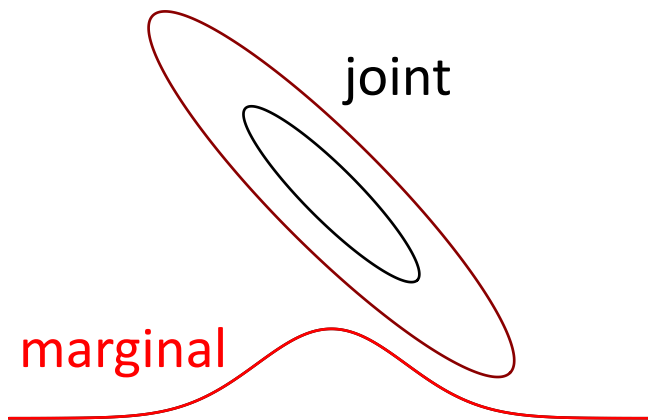
Cov matrix

$$K = \begin{bmatrix} \mathbb{E}[(Y_1 - \mu_1)(Y_1 - \mu_1)] & \cdots & \mathbb{E}[(Y_1 - \mu_1)(Y_n - \mu_n)] \\ \mathbb{E}[(Y_2 - \mu_2)(Y_1 - \mu_1)] & \cdots & \mathbb{E}[(Y_2 - \mu_2)(Y_n - \mu_n)] \\ \vdots & \ddots & \vdots \\ \mathbb{E}[(Y_n - \mu_n)(Y_1 - \mu_1)] & \cdots & \mathbb{E}[(Y_n - \mu_n)(Y_n - \mu_n)] \end{bmatrix}$$

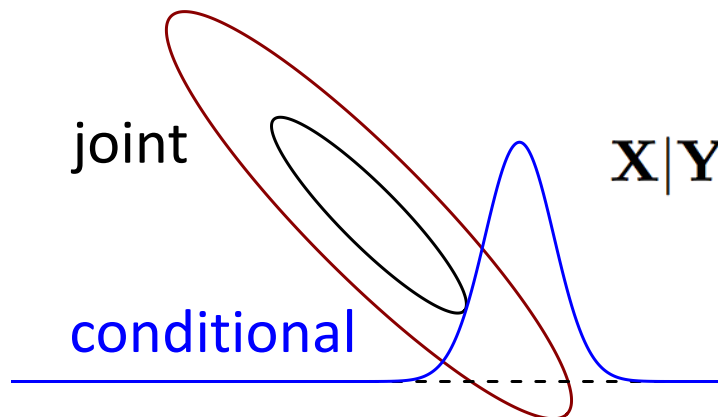
# Motivation: the basics III

$$\mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}, \mathbf{Y}}, \boldsymbol{\Sigma}_{\mathbf{X}, \mathbf{Y}})$$

$$\boldsymbol{\mu}_{\mathbf{X}, \mathbf{Y}} = \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \quad \boldsymbol{\Sigma}_{\mathbf{X}, \mathbf{Y}} = \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}$$



$$\mathbf{X} \sim \mathcal{N}(\mathbf{a}, A) \quad \mathbf{Y} \sim \mathcal{N}(\mathbf{b}, C)$$



$$\mathbf{X} | \mathbf{Y} \sim \mathcal{N}(\mathbf{a} + BC^{-1}(\mathbf{Y} - \mathbf{b}), A - BC^{-1}B^\top)$$

## Motivation: the problem?

---

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \right)$$



## Motivation: the problem?

---

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \right)$$

Joint probability -> conditional probability -> posterior from the prior,  
but here we have the joint probability of the values of  $f(x)$  for all  $x$

“observed”, so: (observed ->  $f$ , non-observed ->  $f_*$ )

# Motivation: the problem?

---

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \right)$$

Joint probability  $\rightarrow$  conditional probability  $\rightarrow$  posterior from the prior, but here we have the joint probability of the values of  $f(x)$  for all  $x$

“observed”, so: (observed  $\rightarrow$   $f$ , non-observed  $\rightarrow$   $f_*$ )

kernel function applied to  $x$   
“similarity” of each  $x_i$  to all  $x_i$

“similarity” of each  $x_i$  to non-observed  $x_i$  that we want to predict

$$\begin{pmatrix} f \\ f_* \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} \Sigma & \Sigma_* \\ \Sigma_*^T & \Sigma_{**} \end{pmatrix} \right)$$

“similarity” of each non-observed  $x_i$  to the other non-observed  $x_i$  that we want to predict

# Motivation: the problem?

---

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \right)$$

Joint probability -> conditional probability -> posterior from the prior,  
but here we have the joint probability of the values of  $f(x)$  for all  $x$

“observed”, so. (observed ->  $f$ , non-observed ->  $f_*$ )  
The magic? (algebra)

From the joint distribution to the conditional one of each  $x_i$  to  
kernel function applied to  $x$  non-observed  $x_i$  that we  
“similarity” of each  $x_i$  to all  $x_i$  want to predict

$$p(f, f_*) \longrightarrow p(f_* | f)$$

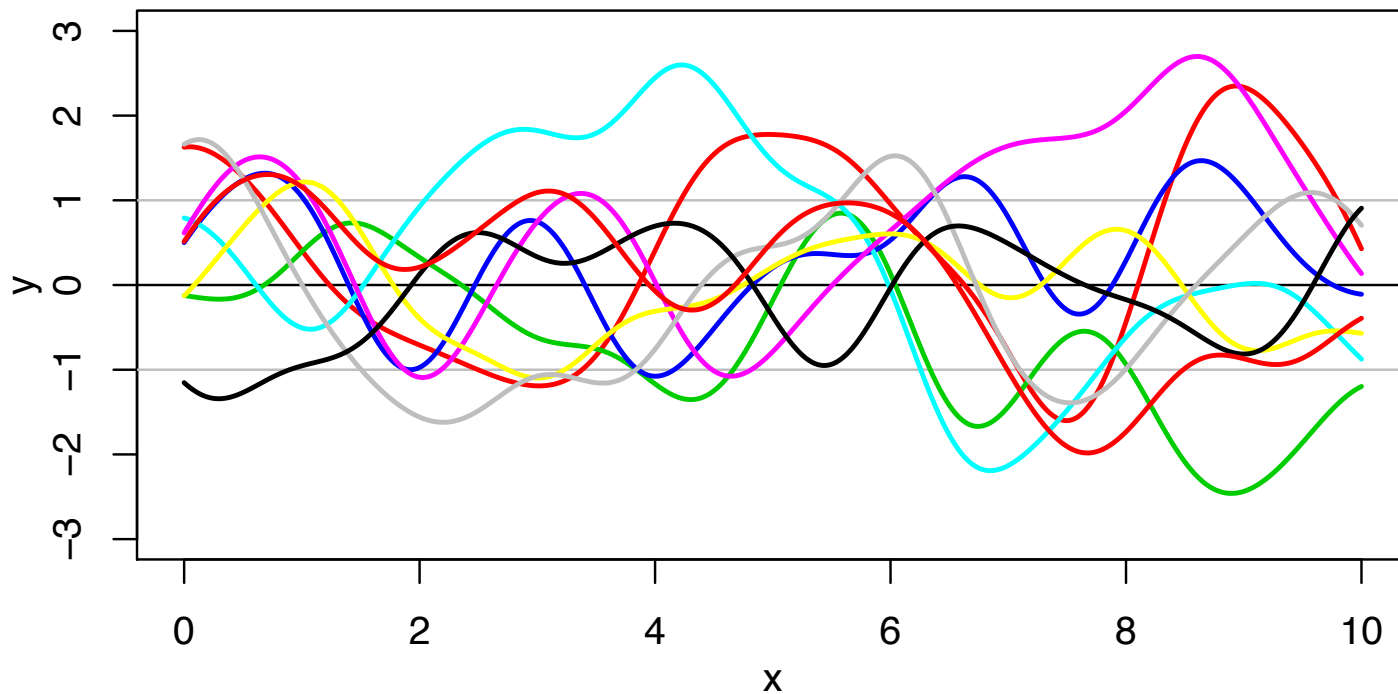
“similarity” of each non-observed  $x_i$   
to the other non-observed  $x_i$  that we  
want to predict

# The definition

---

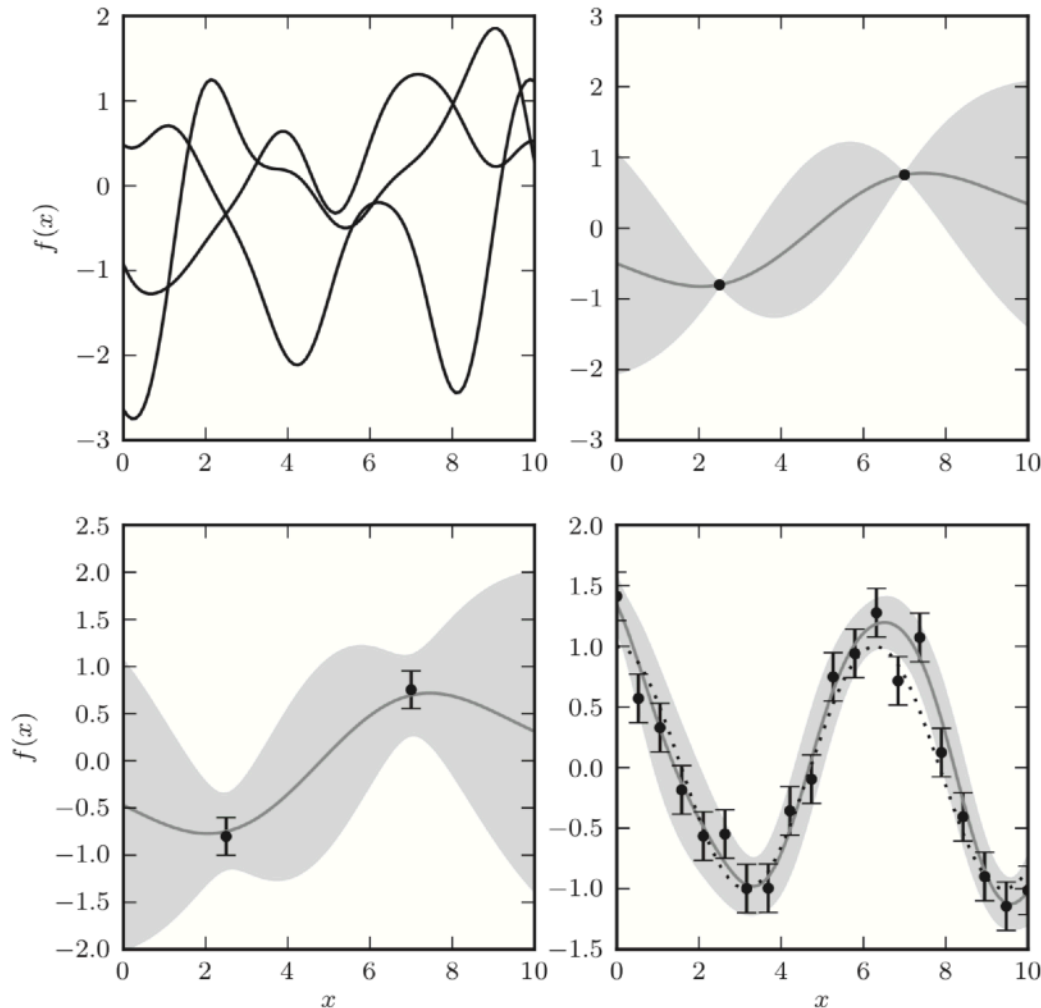
“A Gaussian process is a collection of random variables with the property that the joint distribution of any finite subset is a Gaussian. “

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



# The Interpretation

The Gaussian kernel can be interpreted as a prior on the functions (instead of on the parameters when doing parametric statistics)

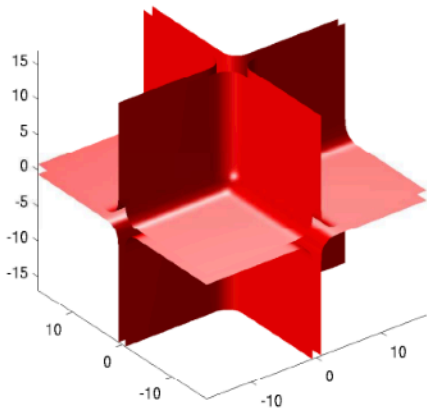


**covariance matrix** ensures that values that are close together in input space will produce output values that are close together.

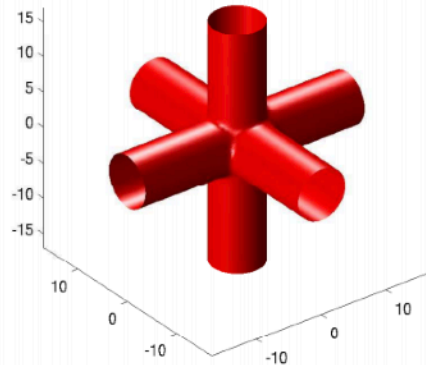
# Is it all about the kernel?

different kernels?!  
additive kernels?!

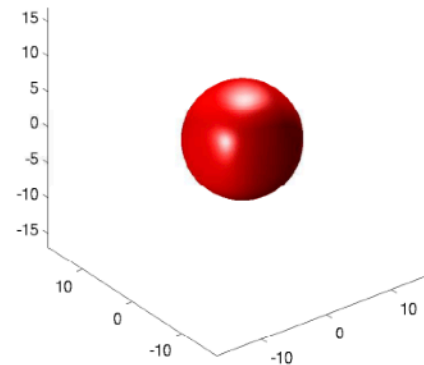
Kind of defining the  
“radius of influence”  
differently at different  
spatial scales



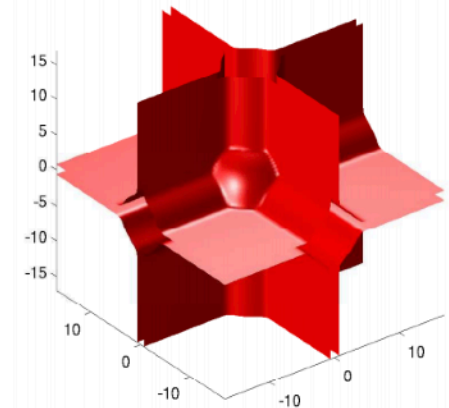
1st order interactions  
 $k_1 + k_2 + k_3$



2nd order interactions  
 $k_1k_2 + k_2k_3 + k_1k_3$



3rd order interactions  
 $k_1k_2k_3$   
(Squared-exp kernel)



All interactions  
(Additive kernel)

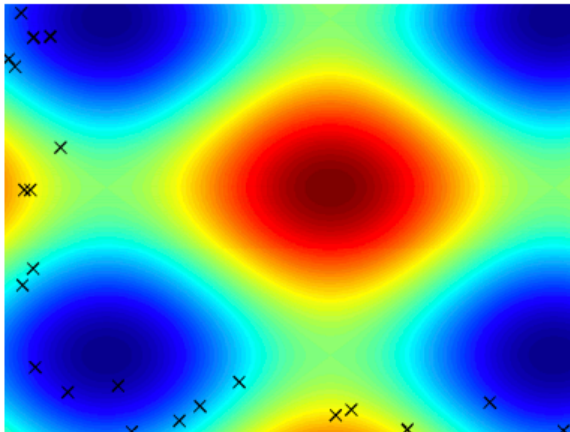
# Is it all about the kernel?

---

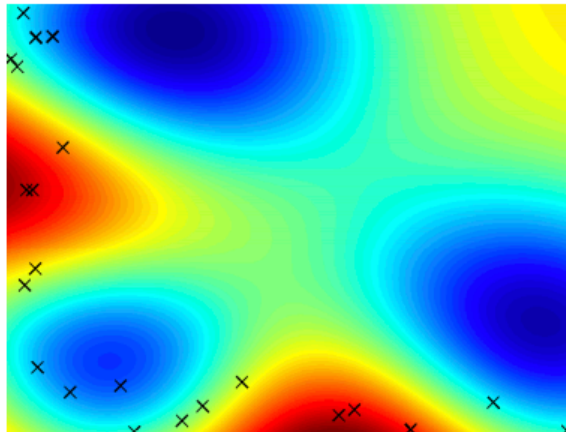
different kernels?!  
additive kernels?!

Kind of defining the  
“radius of influence”  
differently at different  
spatial scales

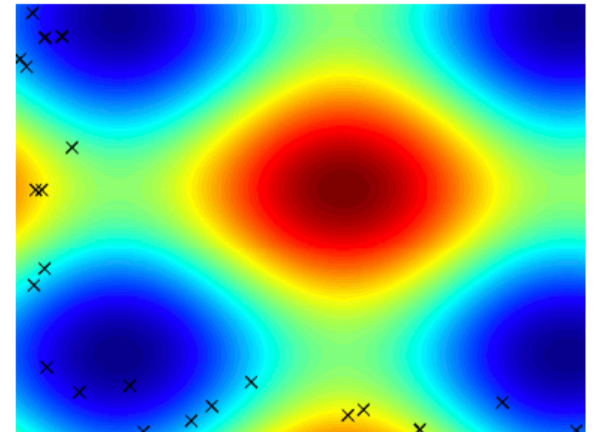
Good at discovering non-local structure!!



True Function  
& data locations



Squared-exp GP  
posterior mean



Additive GP  
posterior mean

# The implementation?

---