

# Stochastic Processes: Lecture 3B

Andrew J. Holbrook

UCLA Biostatistics 270

Spring 2022

# Defining GPs

## Definition 1

A r.v.  $X$  on a Hilbert space  $H$  is Gaussian, if for any  $h \in H$  the r.v.  $\langle X, h \rangle$  follows a Gaussian distribution on  $\mathbb{R}$ .

## Proposition 1

If  $X$  is a r.v. on  $H$ , then there exist an  $m \in H$  and a positive, symmetric, nuclear (trace-class) operator  $K$  on  $H$  such that:

$$E(\langle X, h \rangle) = \langle m, h \rangle, \quad h \in H,$$

$$E(\langle X - m, h \rangle \langle X - m, g \rangle) = \langle h, Kg \rangle, \quad h, g \in H.$$

$m$  and  $K$  are uniquely determined. Conversely, for  $m$  and  $K$  defined as above, there exists a unique Gaussian distribution on  $H$  satisfying these moment conditions.

# A definition suited for GP regression

## Definition 2

*A Gaussian process is a collection of r.v.s, any finite number of which have a joint Gaussian distribution.*

A GP  $f(x)$  is completely specified by its mean and covariance functions,  $m(x)$  and  $K(x, x')$ . These satisfy

$$\begin{aligned}m(x) &= E(f(x)), \\K(x, x') &= E((f(x) - m(x))(f(x') - m(x')))).\end{aligned}$$

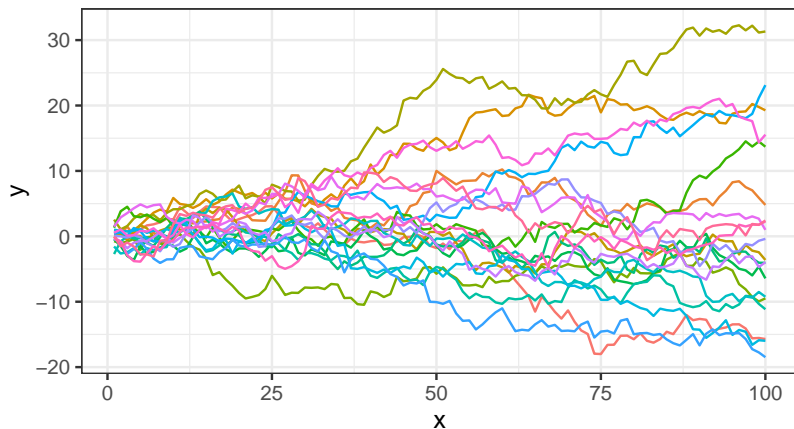
It is common to use the notation

$$f(x) \sim \mathcal{GP}(m(x), K(x, x')) .$$

# Common Kernels

Brownian motion kernel  $K(t, s) = \min(t, s)$ .

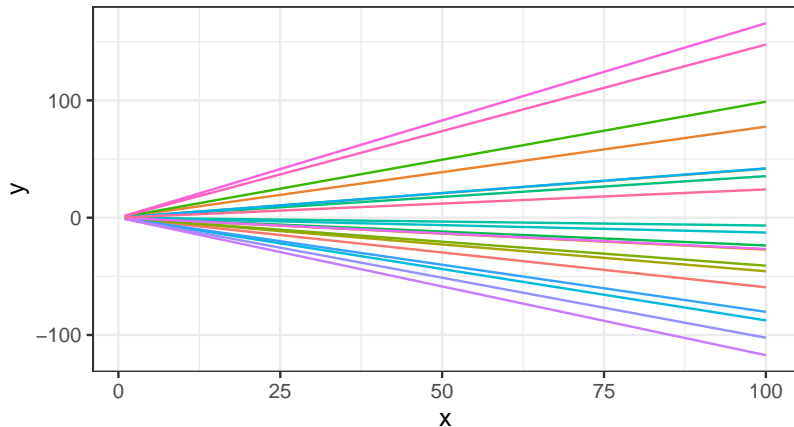
Brownian Motions



# Common Kernels

Linear kernel  $K(t, s) = ts$  or  $K = xx^T$  for vector  $x$ .

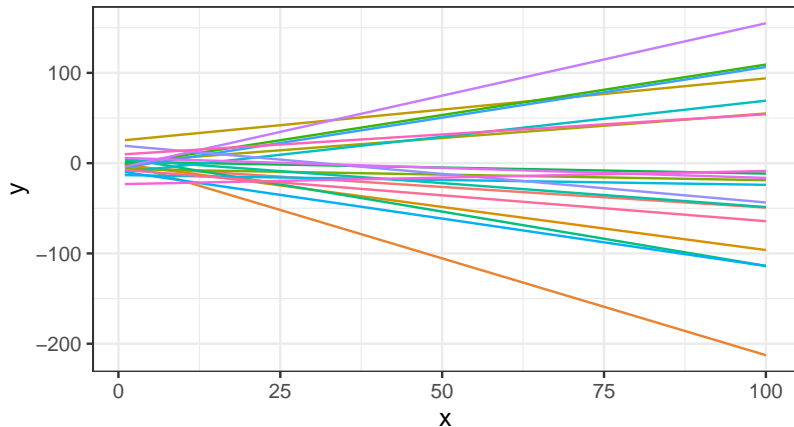
## Linear Kernel



# Common Kernels

Linear kernel  $K = XX^T$  for  $X$  a design matrix  $[\alpha \mathbf{1}, \mathbf{x}]$ .

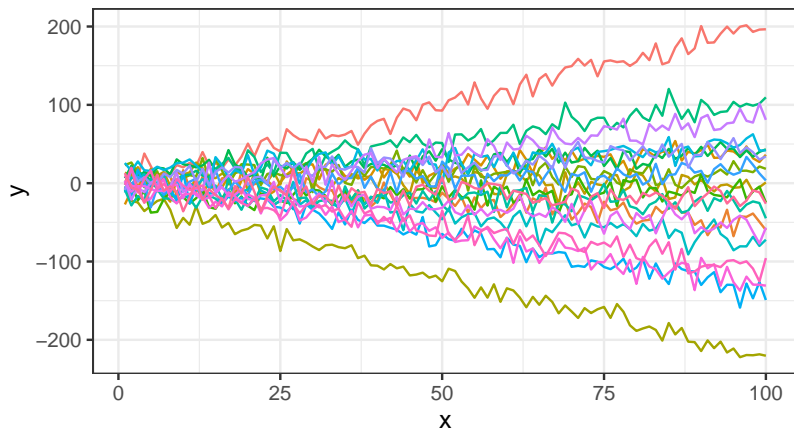
## Linear Kernel with Intercept



# Common Kernels

Linear kernel  $K = XX^T + \sigma^2 I$  for  $X$  a design matrix  $[\alpha \mathbf{1}, \mathbf{x}]$ .

## Linear Kernel with Intercept and Noise Term



# Common Kernels

Matérn kernel takes  $d = d(x, x')$  the distance between  $x$  and  $x'$ :

$$K_\nu(d) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \sqrt{2\nu} \frac{d}{\rho} \right)^\nu C_\nu \left( \sqrt{2\nu} \frac{d}{\rho} \right).$$

Here,  $\sigma^2$  is the variance term,  $C_\nu$  is the modified Bessel function of the second kind, and  $\rho$  is the lengthscale. A GP with covariance  $K_\nu$  is  $\lceil \nu \rceil - 1$  times differentiable.

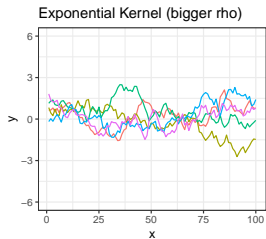
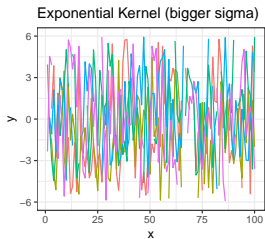
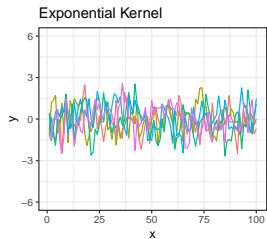
Examples:

- ▶  $\nu = 1/2$ :  $K_\nu(d) = \sigma^2 e^{-d/\rho}$  (exponential kernel)
- ▶  $\nu = 3/2$ :  $K_\nu(d) = \sigma^2 (1 + \sqrt{3}d/\rho) e^{-\sqrt{3}d/\rho}$
- ▶  $\lim_{\nu \rightarrow \infty} K_\nu(d) = \sigma^2 e^{-d^2/(2\rho^2)}$   
(squared exponential or radial basis function kernel)



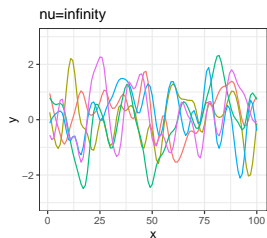
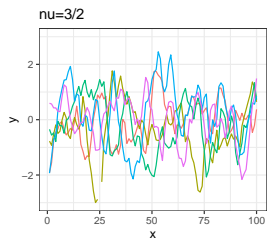
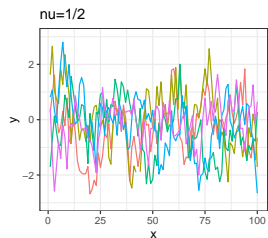
# Common Kernels

Exponential kernel  $K_{1/2}(d) = \sigma^2 e^{-d/\rho}$ .



# Common Kernels

Matérn kernels:  $\nu = 1/2$ ,  $\nu = 3/2$ ,  $\nu = \infty$ .



# Building Kernels

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

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbf{x}$  to  $\mathbb{R}^M$ ,  $k_3(\cdot, \cdot)$  is a valid kernel in  $\mathbb{R}^M$ ,  $\mathbf{A}$  is a symmetric positive semidefinite matrix,  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are variables (not necessarily disjoint) with  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a$  and  $k_b$  are valid kernel functions over their respective spaces.

- Bishop (2006). “Pattern Recognition and Machine Learning”.

## Prediction

Prediction follows from the conditional distribution of a multivariate Gaussian. Write  $K(x, x) := \text{cov}(f(x), f(x))$  (an abuse of notation) and suppose we have observed

$$f(x) \sim N_D(0, K(x, x))$$

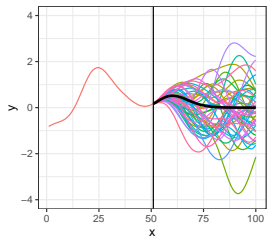
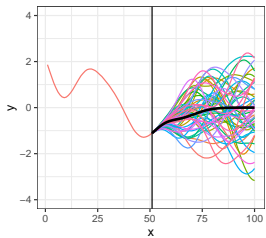
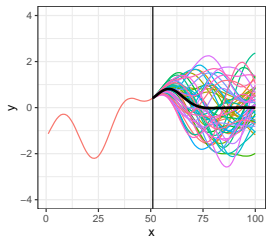
and we would like to predict  $f(x^*) \sim N_d(0, K(x^*, x^*))$ . Then we use the joint distribution

$$\begin{pmatrix} f(x) \\ f(x^*) \end{pmatrix} = N_{D+d} \left( 0, \begin{pmatrix} K(x, x) & K(x, x^*) \\ K(x^*, x) & K(x^*, x^*) \end{pmatrix} \right)$$

to conclude that

$$f(x^*)|f(x) \sim N_d \left( K(x^*, x)K(x, x)^{-1}f(x), \right. \\ \left. K(x^*, x^*) - K(x^*, x)K(x, x)^{-1}K(x, x^*) \right).$$

# Prediction



# Bayesian inference

Suppose we've observed  $N$  pairs  $(f, X) = [(f_1, x_1), \dots, (f_N, x_N)]$ .  
Denote  $\theta = (\sigma^2, \rho, \tau^2)$  for the kernel

$$K_\theta(x_n, x_{n'}) = \sigma^2 e^{-\frac{1}{2\rho^2}(x_n - x_{n'})^2} + \tau^2 \delta_{nn'}.$$

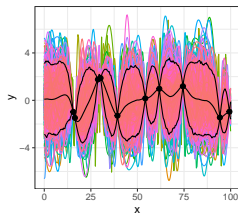
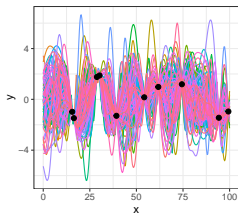
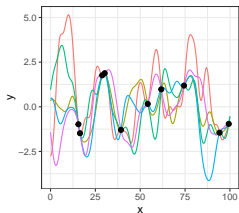
We know  $f|X, \theta \sim N_N(0, K_\theta(X))$ , where  $K_\theta(X) = [K_\theta(x_n, x_{n'})]_{n,n'=1}^N$   
is the  $N \times N$  covariance matrix. The likelihood is

$$p(f|X, \theta) \propto |K_\theta(X)|^{-1/2} e^{-\frac{1}{2}f^T K_\theta^{-1}(X)f}.$$

Specify priors  $p(\sigma^2)$ ,  $p(\rho)$ ,  $p(\tau^2)$  and the posterior becomes

$$p(\theta|f, X) \propto |K_\theta(X)|^{-1/2} e^{-\frac{1}{2}f^T K_\theta^{-1}(X)f} p(\sigma^2)p(\rho)p(\tau^2).$$

# Posterior predictive curves, mean and intervals



## Binary classification

Now suppose we've observed  $N$  pairs

$(y, \mathbf{X}) = [(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)]$  for  $y_i$  binary. We posit Gaussian latent variables  $f_1, \dots, f_n$ . Use MCMC to infer

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{f} | y, \mathbf{X}) &\propto |K_{\boldsymbol{\theta}}(\mathbf{X})|^{-1/2} e^{-\frac{1}{2} \mathbf{f}^T K_{\boldsymbol{\theta}}^{-1}(\mathbf{X}) \mathbf{f}} \\ &\quad + p(\sigma^2) p(\rho) p(\tau^2) \\ &\quad + \prod_i \left( \frac{e^{f_i}}{1 + e^{f_i}} \right)^{y_i} \left( \frac{1}{1 + e^{f_i}} \right)^{1-y_i} \end{aligned}$$

Note the log of last line simplifies to

$$\sum_i \left( y_i f_i - \log(1 + e^{f_i}) \right) .$$

One conditions on posterior samples of  $\boldsymbol{\theta}$  and  $\mathbf{f}$  to get posterior predictive curves  $f(x^*)$  for  $x^*$  unobserved.



# Binary classification

Posterior predictive curves are no longer “pinned” down by observations, since these themselves are inferred latent variables.

