

## Spectral kernels

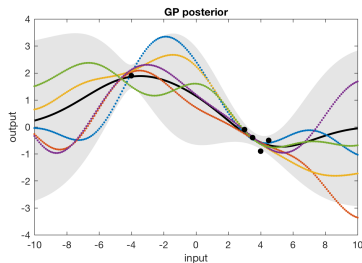
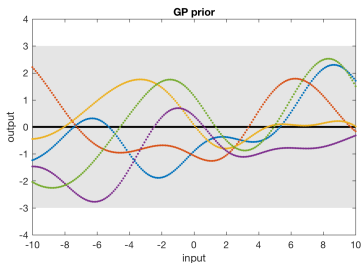
Markus Heinonen, Research fellow  
Aalto University, Finland

GPSS'21  
13.9.2021

# Gaussian processes

- Bayesian non-parametric kernel model
- Key idea: **function prior**  $f(x) \sim \mathcal{GP}(m(x), K_\theta(x, x'))$  that encodes

$$p \left( \begin{matrix} f(x_1) \\ \vdots \\ f(x_N) \end{matrix} \right) = \mathcal{N} \left( \underbrace{\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{bmatrix}}_{\mathbf{f}} \mid \underbrace{\begin{bmatrix} m(x_1) \\ \vdots \\ m(x_N) \end{bmatrix}}_{\mathbf{m}}, \underbrace{\begin{bmatrix} K_\theta(x_1, x_1) & \cdots & K_\theta(x_1, x_N) \\ \vdots & \ddots & \vdots \\ K_\theta(x_N, x_1) & \cdots & K_\theta(x_N, x_N) \end{bmatrix}}_{K_\theta} \right)$$



## How to learn a kernel?

- Choose prior with maximum **volume** of data-matching functions

$$\underbrace{\log p(\mathbf{y}|\theta)}_{\text{marginal log likelihood}} = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (1)$$

$$= -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi \quad (2)$$

- Relatively robust against overfitting
  - Determinant: volume of space spanned by kernel
  - Finds a simple basis for the data
  - Overfitting still possible, if  $p(\mathbf{f})$  can be shaped to match  $p(\mathbf{y}|\mathbf{f})$
- Powerful formalism to learn kernels
  - Replaces cross-validation
  - We can (auto)differentiate  $\log p(\mathbf{y}|\theta)$  and optimise wrt  $\theta$
  - GPflow, GPyTorch, Stan, GaussianProcesses.jl, etc

## How to learn a kernel?

- Choose prior with maximum **volume** of data-matching functions

$$\underbrace{\log p(\mathbf{y}|\theta)}_{\text{marginal log likelihood}} = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (1)$$

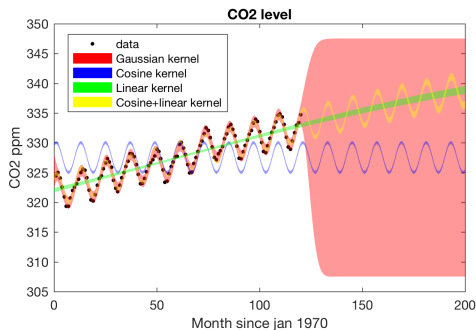
$$= -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi \quad (2)$$

- Relatively robust against overfitting
  - Determinant: volume of space spanned by kernel
  - Finds a simple basis for the data
  - Overfitting still possible, if  $p(\mathbf{f})$  can be shaped to match  $p(\mathbf{y}|\mathbf{f})$
- Powerful formalism to learn kernels
  - Replaces cross-validation
  - We can (auto)differentiate  $\log p(\mathbf{y}|\theta)$  and optimise wrt  $\theta$
  - GPflow, GPyTorch, Stan, GaussianProcesses.jl, etc



## How to choose kernel?

- Gaussian kernel  $K_g(x, x') = \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$
- Periodic kernel  $K_{cos}(x, x') = \exp\left(-\frac{2 \sin^2(\pi|x-x'|/p)}{\ell^2}\right)$
- Linear kernel  $K_{lin}(x, x') = xx' + c$
- Composite kernel, eg  $K(x, x') = K_g(x, x') + K_{lin}(x, x')$



- Our topic: **Spectral kernels** can learn **arbitrary** kernel function forms

# Fourier transforms

- **Fourier transform**  $S(\omega)$  of a function  $f(x)$ ,

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i x \omega} dx \quad (3)$$

where

- $i$  is the imaginary number with  $i^2 = -1$  and  $i^0 = 1$
- $\omega$  is a frequency
- **Inverse Fourier transform**  $f(x)$  of spectral density  $S(\omega)$ ,

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi i x \omega} d\omega \quad (4)$$

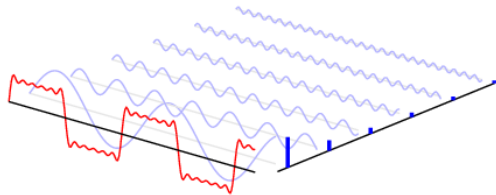
- **Euler's identity** helps compute Fourier's in practise

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + \underbrace{i \sin x}_{\text{imaginary part}} \quad (5)$$

where the complex part often cancels out

- Hence,

$$e^{\pm 2\pi i x \omega} = \cos(2\pi x \omega) \pm i \sin(2\pi x \omega) \quad (6)$$



# Fourier transforms

- **Fourier transform**  $S(\omega)$  of a function  $f(x)$ ,

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i x \omega} dx \quad (3)$$

where

- $i$  is the imaginary number with  $i^2 = -1$  and  $i^0 = 1$
- $\omega$  is a frequency
- **Inverse Fourier transform**  $f(x)$  of spectral density  $S(\omega)$ ,

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi i x \omega} d\omega \quad (4)$$

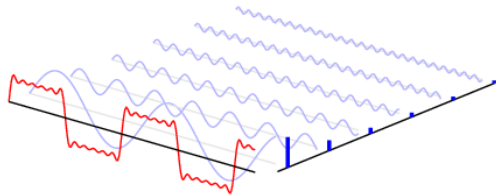
- **Euler's identity** helps compute Fourier's in practise

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + \underbrace{i \sin x}_{\text{imaginary part}} \quad (5)$$

where the complex part often cancels out

- Hence,

$$e^{\pm 2\pi i x \omega} = \cos(2\pi x \omega) \pm i \sin(2\pi x \omega) \quad (6)$$



## Fourier duals for kernels

- Let's apply Fourier to the kernel  $K(\tau) := K(x, x')$ , where  $\tau = x - x'$  (instead of  $f(x)$ )

### Theorem (Bochner)

Any *stationary kernel*  $K : \mathbb{R}^D \mapsto \mathbb{R}$  and its *spectral density*  $S : \mathbb{R}^D \mapsto \mathbb{R}$  are Fourier duals

$$K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

$$S(\omega) = \int_{-\infty}^{\infty} K(\tau) e^{-2\pi i \omega^T \tau} d\tau. \quad (\text{Fourier Transform})$$

- 1 All stationary kernels have *spectral density*  $S(\omega)$ 
  - If someone gives you a kernel  $K(\tau)$ , we can solve what frequencies it considers by solving the FT
  - Spectral features are of theoretical interest
- 2 All spectral densities define a covariance function  $K(\tau)$ 
  - If someone gives you a spectral density  $S(\omega)$ , we can solve its similarity function (=kernel) by solving the IFT
  - If we change the spectral density, we get a new kernel
  - $\Rightarrow$  kernel learning

## Fourier duals for kernels

- Let's apply Fourier to the kernel  $K(\tau) := K(x, x')$ , where  $\tau = x - x'$  (instead of  $f(x)$ )

### Theorem (Bochner)

Any *stationary kernel*  $K : \mathbb{R}^D \mapsto \mathbb{R}$  and its *spectral density*  $S : \mathbb{R}^D \mapsto \mathbb{R}$  are Fourier duals

$$K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

$$S(\omega) = \int_{-\infty}^{\infty} K(\tau) e^{-2\pi i \omega^T \tau} d\tau. \quad (\text{Fourier Transform})$$

- 1 All stationary kernels have *spectral density*  $S(\omega)$ 
  - If someone gives you a kernel  $K(\tau)$ , we can solve what frequencies it considers by solving the FT
  - Spectral features are of theoretical interest
- 2 All spectral densities define a covariance function  $K(\tau)$ 
  - If someone gives you a spectral density  $S(\omega)$ , we can solve its similarity function (=kernel) by solving the IFT
  - If we change the spectral density, we get a new kernel
  - $\Rightarrow$  kernel learning

## Fourier duals for kernels

- Let's apply Fourier to the kernel  $K(\tau) := K(x, x')$ , where  $\tau = x - x'$  (instead of  $f(x)$ )

### Theorem (Bochner)

Any *stationary kernel*  $K : \mathbb{R}^D \mapsto \mathbb{R}$  and its *spectral density*  $S : \mathbb{R}^D \mapsto \mathbb{R}$  are Fourier duals

$$K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

$$S(\omega) = \int_{-\infty}^{\infty} K(\tau) e^{-2\pi i \omega^T \tau} d\tau. \quad (\text{Fourier Transform})$$

- 1 All stationary kernels have *spectral density*  $S(\omega)$ 
  - If someone gives you a kernel  $K(\tau)$ , we can solve what frequencies it considers by solving the FT
  - Spectral features are of theoretical interest
- 2 All spectral densities define a covariance function  $K(\tau)$ 
  - If someone gives you a spectral density  $S(\omega)$ , we can solve its similarity function (=kernel) by solving the IFT
  - If we change the spectral density, we get a new kernel
  - $\Rightarrow$  kernel learning

## Fourier duals for kernels

- Let's apply Fourier to the kernel  $K(\tau) := K(x, x')$ , where  $\tau = x - x'$  (instead of  $f(x)$ )

### Theorem (Bochner)

Any *stationary kernel*  $K : \mathbb{R}^D \mapsto \mathbb{R}$  and its *spectral density*  $S : \mathbb{R}^D \mapsto \mathbb{R}$  are Fourier duals

$$K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

$$S(\omega) = \int_{-\infty}^{\infty} K(\tau) e^{-2\pi i \omega^T \tau} d\tau. \quad (\text{Fourier Transform})$$

- 1 All stationary kernels have *spectral density*  $S(\omega)$ 
  - If someone gives you a kernel  $K(\tau)$ , we can solve what frequencies it considers by solving the FT
  - Spectral features are of theoretical interest
- 2 All spectral densities define a covariance function  $K(\tau)$ 
  - If someone gives you a spectral density  $S(\omega)$ , we can solve its similarity function (=kernel) by solving the IFT
  - If we change the spectral density, we get a new kernel
  - $\Rightarrow$  kernel learning

## Kernel sinusoid representation

- Assume symmetric frequency distribution  $S(\omega) = S(-\omega)$
- Euler's identity  $e^{\pm ix} = \cos x \pm i \sin x$
- Sine identity  $\sin(-x) = -\sin(x)$
- Then we can solve the inverse Fourier as

$$\begin{aligned}K(\tau) &= \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \tau \omega} d\omega \\&= \int_{-\infty}^{\infty} S(\omega) \cos(2\pi \tau \omega) d\omega + \int_{-\infty}^{\infty} i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_{-\infty}^0 i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega + \int_0^{\infty} i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_0^{\infty} i S(-\omega) \sin(2\pi \tau (-\omega)) d\omega + \int_0^{\infty} i S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_0^{\infty} -i S(\omega) \sin(2\pi \tau \omega) d\omega + \int_0^{\infty} i S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega)\end{aligned}$$

- All stationary kernels are linear  $S(\omega)$ -combinations of  $\cos(2\pi \tau \omega)$



## Kernel sinusoid representation

- Assume symmetric frequency distribution  $S(\omega) = S(-\omega)$
- Euler's identity  $e^{\pm ix} = \cos x \pm i \sin x$
- Sine identity  $\sin(-x) = -\sin(x)$
- Then we can solve the inverse Fourier as

$$\begin{aligned}K(\tau) &= \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \tau \omega} d\omega \\&= \int_{-\infty}^{\infty} S(\omega) \cos(2\pi \tau \omega) d\omega + \int_{-\infty}^{\infty} i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_{-\infty}^0 i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega + \int_0^{\infty} i \cdot S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_0^{\infty} i S(-\omega) \sin(2\pi \tau (-\omega)) d\omega + \int_0^{\infty} i S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega) + \int_0^{\infty} -i S(\omega) \sin(2\pi \tau \omega) d\omega + \int_0^{\infty} i S(\omega) \sin(2\pi \tau \omega) d\omega \\&= \mathbb{E}_{S(\omega)} \cos(2\pi \tau \omega)\end{aligned}$$

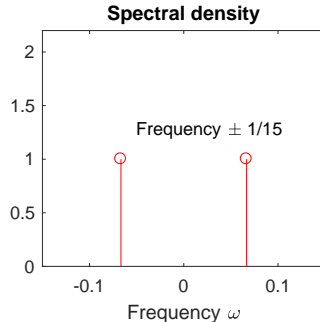
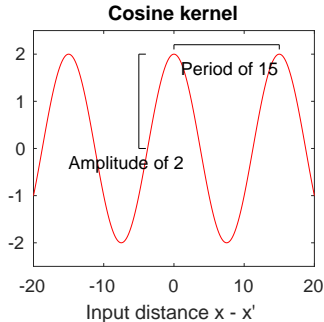
- All stationary kernels are linear  $S(\omega)$ -combinations of  $\cos(2\pi \tau \omega)$

## Kernel sinusoid representation

- Our new general **kernel definition**

$$K(\tau) = \mathbb{E}_{S(\omega)} \cos(2\pi\tau\omega) \quad (7)$$

- Frequency  $\omega$  is inverse of period  $1/\omega$
- Amplitude  $S(\omega)$
- Frequencies are symmetric  $S(\omega) = S(-\omega)$
- With  $S(\omega) = \delta_{1/15}(\omega)$ , the kernel becomes  $K(\tau) = \cos(2\pi\tau \frac{1}{15})$



## Gaussian kernel sinusoids

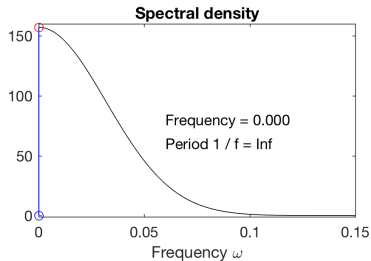
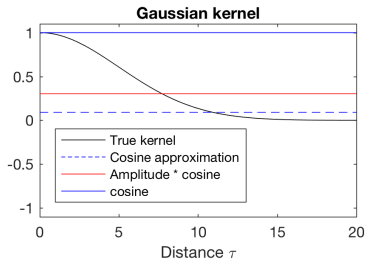
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (8)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (9)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (10)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (11)$$



## Gaussian kernel sinusoids

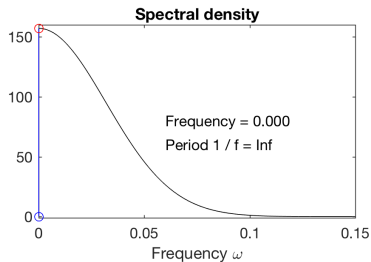
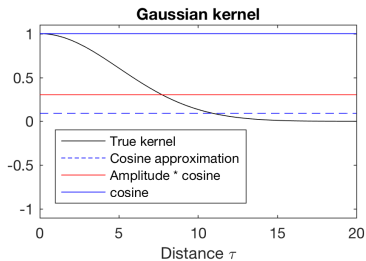
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

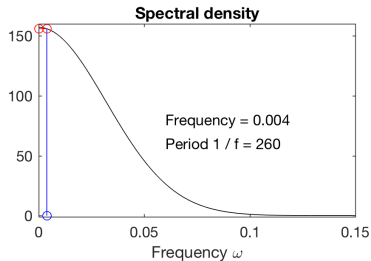
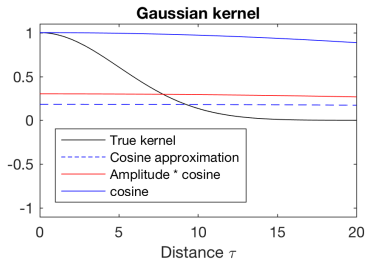
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

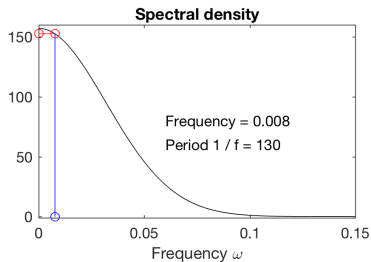
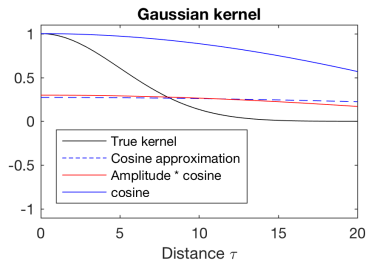
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

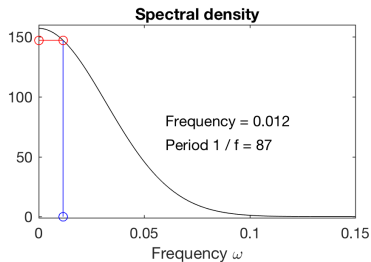
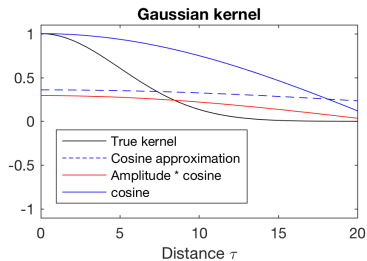
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

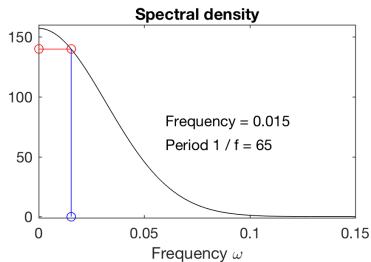
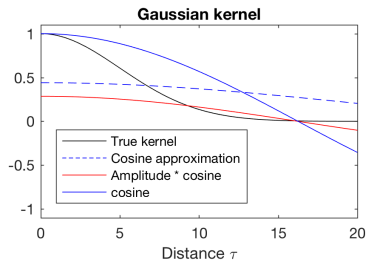
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$





## Gaussian kernel sinusoids

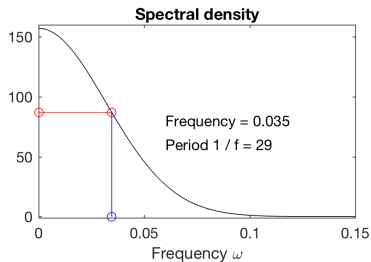
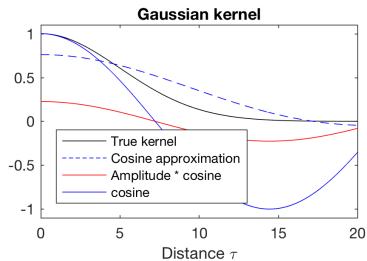
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

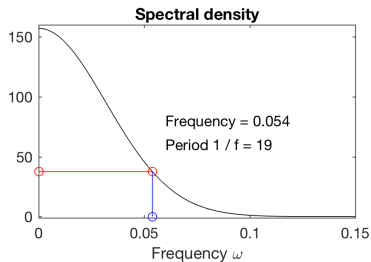
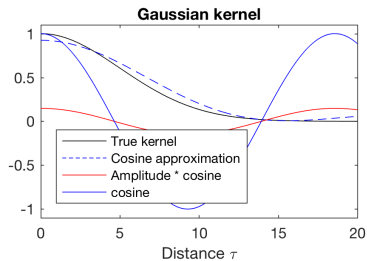
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

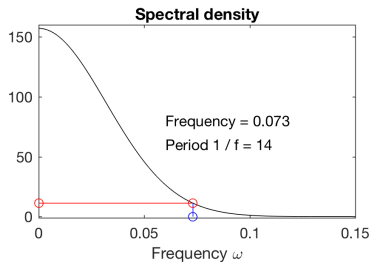
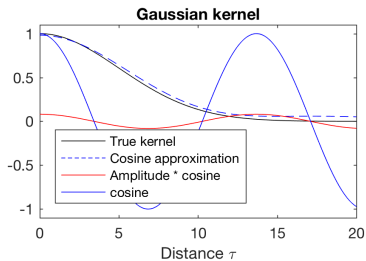
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

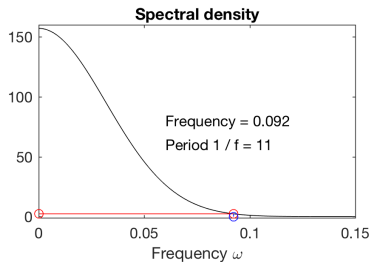
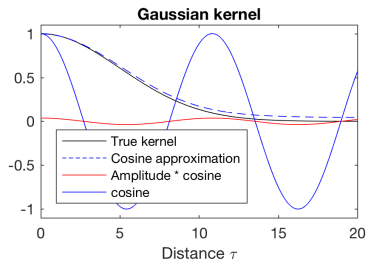
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

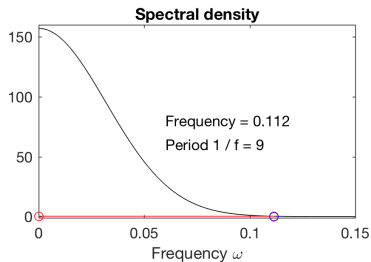
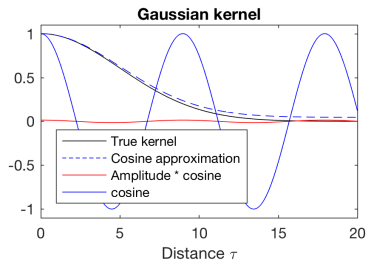
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

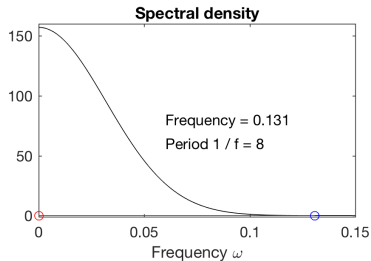
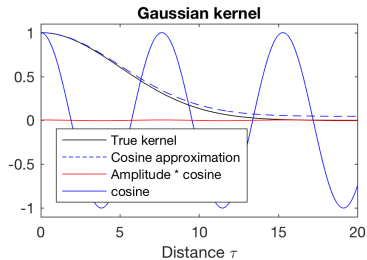
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Gaussian kernel sinusoids

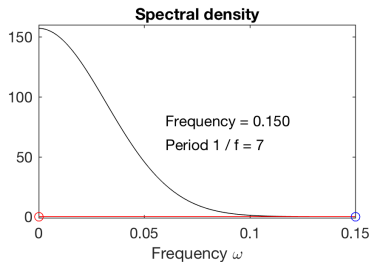
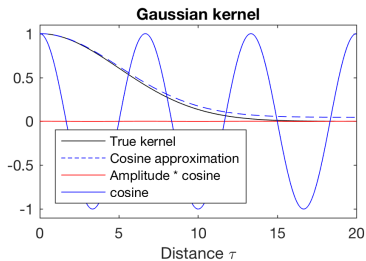
- Gaussian kernel  $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$  fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (12)$$

$$= 2\pi \ell^2 \exp(-2\pi^2 \ell^2 \omega^2) \quad (13)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (14)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (15)$$



## Some spectral densities

$$K_{gauss}(\tau) = \exp\left(-\frac{\tau^2}{\ell^2}\right)$$

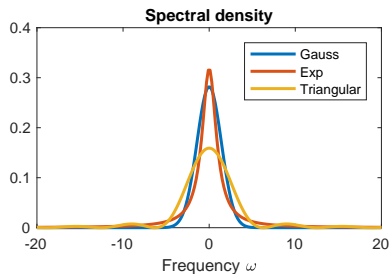
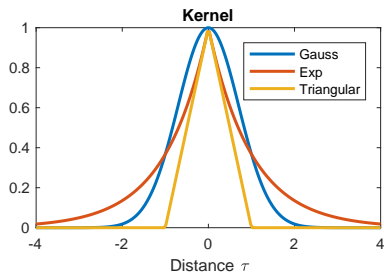
$$K_{exp}(\tau) = \exp(-|\tau|/\ell)$$

$$K_{tri}(\tau) = 0.5(1 - |\tau|)_+$$

$$S_{gauss}(\omega) = \frac{\sqrt{\ell}}{2\sqrt{\pi}} \exp(-\ell\omega^2/4) \quad (16)$$

$$S_{exp}(\omega) = 1/(\pi/\ell + \pi\ell\omega^2) \quad (17)$$

$$S_{tri}(\omega) = (1 - \cos \omega)/(\pi\omega^2) \quad (18)$$



- Can we construct **new** kernels from custom spectral densities?



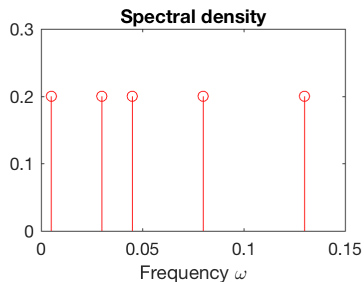
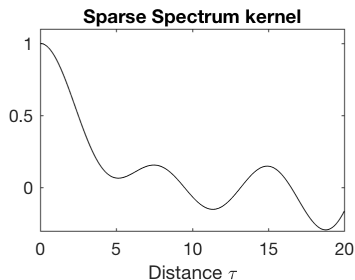
## Sparse Spectrum (SS) kernel<sup>1</sup>

- Define  $Q$  real frequencies  $(\omega_1, \dots, \omega_Q)^T \in \mathbb{R}^Q$  with Fourier dual

$$S(\omega) := \frac{1}{Q} \sum_{i=1}^Q \delta(\omega - \omega_i) \quad (19)$$

$$\xrightarrow{\text{IFT}} K(\tau) = \frac{1}{Q} \sum_{i=1}^Q \cos(2\pi\tau\omega_i) \quad (20)$$

- No decay, prone to overfitting



<sup>1</sup>Lazaro-Gredilla et al (JMLR 2010) Sparse spectrum gaussian process regression

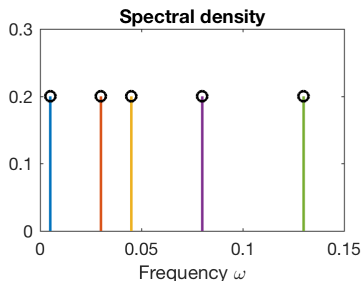
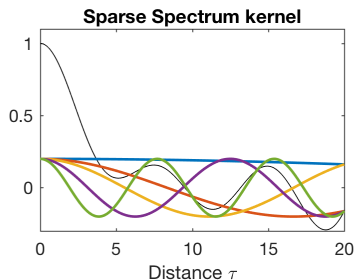
## Sparse Spectrum (SS) kernel<sup>1</sup>

- Define  $Q$  real frequencies  $(\omega_1, \dots, \omega_Q)^T \in \mathbb{R}^Q$  with Fourier dual

$$S(\omega) := \frac{1}{Q} \sum_{i=1}^Q \delta(\omega - \omega_i) \quad (19)$$

$$\xrightarrow{\text{IFT}} K(\tau) = \frac{1}{Q} \sum_{i=1}^Q \cos(2\pi\tau\omega_i) \quad (20)$$

- No decay, prone to overfitting



<sup>1</sup>Lazaro-Gredilla et al (JMLR 2010) Sparse spectrum gaussian process regression

## Spectral Mixture (SM) kernel<sup>2</sup>

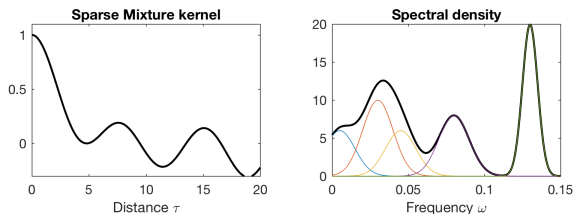
- Define mixture of  $Q$  Gaussians  $\{a_i \mathcal{N}(\mu_i, \sigma_i^2)\}_{i=1}^Q$

$$S(\omega) := \sum_{i=1}^Q a_i \mathcal{N}(\omega | \mu_i, \sigma_i^2) \quad (21)$$

$$\xrightarrow{\text{IFT}} K(\tau) = \int_{-\infty}^{\infty} S(\omega) \cos(2\pi\tau\omega) d\omega \quad (22)$$

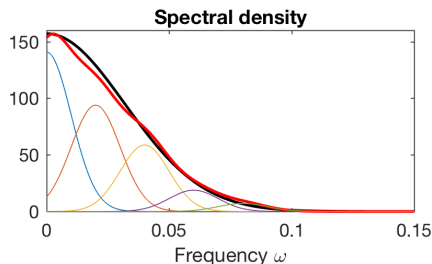
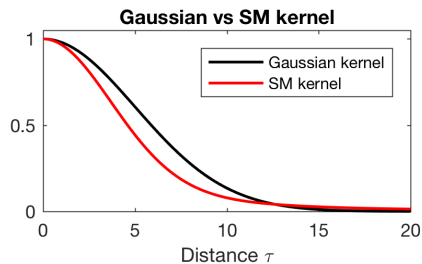
$$= \sum_{i=1}^Q a_i \underbrace{\exp(-2\pi^2 \sigma_i^2 \tau^2)}_{\text{smooth decay}} \underbrace{\cos(2\pi\tau\mu_i)}_{\text{periodic}} \quad (23)$$

- Dense in the set of stationary kernels  $\Rightarrow$  can generate any stationary kernel



<sup>2</sup>Wilson, Adams (ICML 2013) Gaussian process kernels for pattern discovery and extrapolation

## Spectral Mixture (SM) kernel



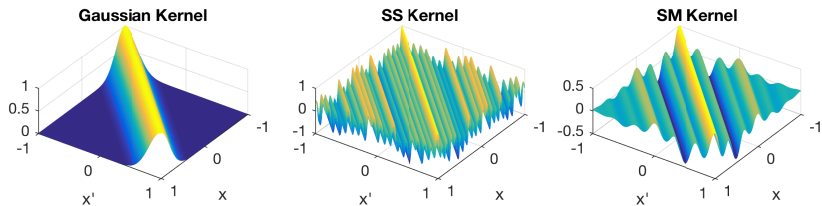
- Approximate gaussian kernel with SM kernel with  $Q = 5$  components, i.e.

$$\sum_{i=1}^Q a_i \exp(-2\pi^2 \sigma_i^2 \tau^2) \cos(2\pi \tau \mu_i) \approx \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$$

for appropriate  $\{a_i, \mu_i, \sigma_i\}$

- (Doable with  $Q = 1$  as well)

## Spectral kernels



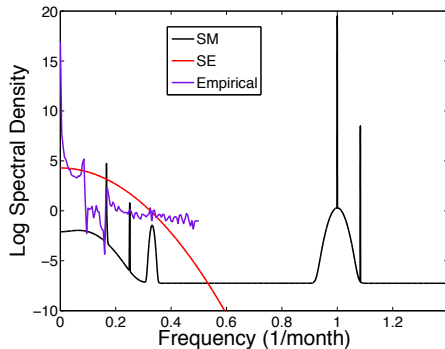
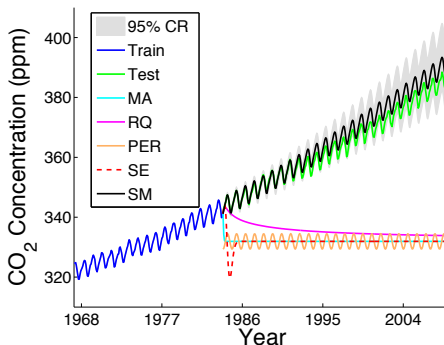
- Image from Remes, Heinonen, Kaski: Non-stationary spectral kernels, NIPS'17

## SM kernel inference

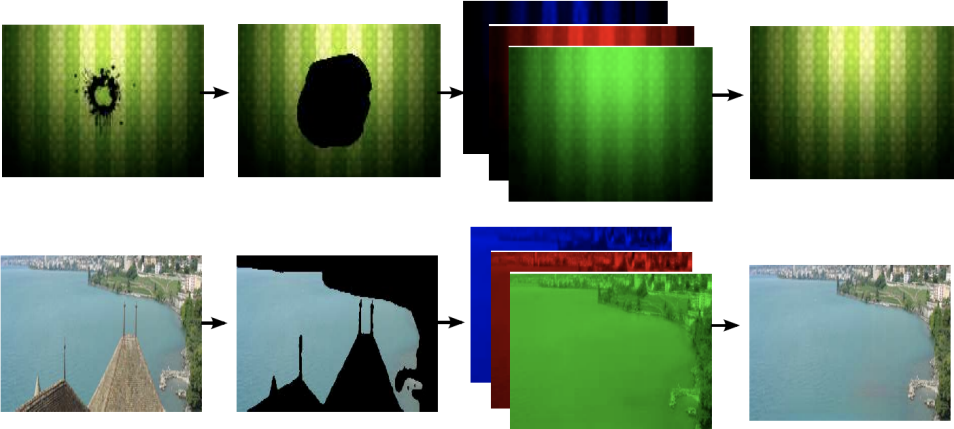
- Optimize  $3Q$  hyperparameters  $\theta = \{a_i, \mu_i, \sigma_i\}_{i=1}^Q$  of kernel  
 $K_\theta(x - x') = \sum_{i=1}^Q a_i \exp(-2\pi^2 \sigma_i^2 \tau^2) \cos(2\pi\tau \mu_i)$  by maximizing

$$\log p(\mathbf{y}|\theta) = -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi$$

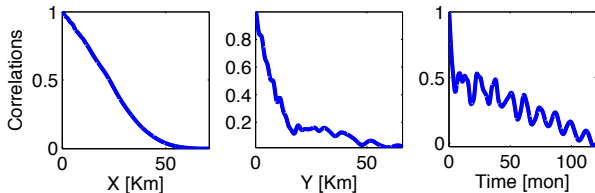
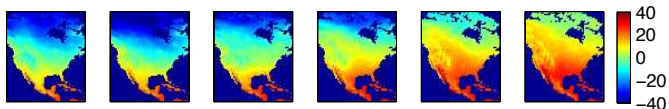
- After kernel is fixed, predictions have closed form



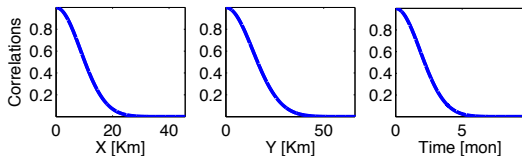
# Image inpainting use case



## Spatio-temporal temperatures



(a) Learned GPatt Kernel for Temperatures



(b) Learned GP-SE Kernel for Temperatures

- SM kernel induces only stationary covariances, but temperatures are non-stationary



## Stationary kernels

- Stationary kernels are **translation-invariant**:

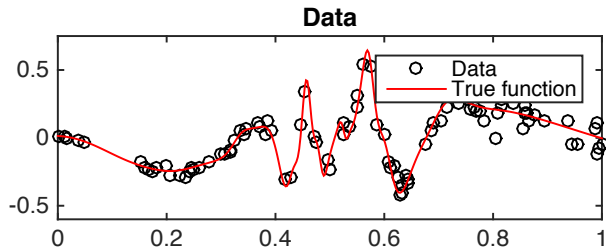
$$K(x, x') = K(x + a, x' + a) \quad (24)$$

$$K(x, x') = K(x - x') \quad (25)$$

for any  $a$

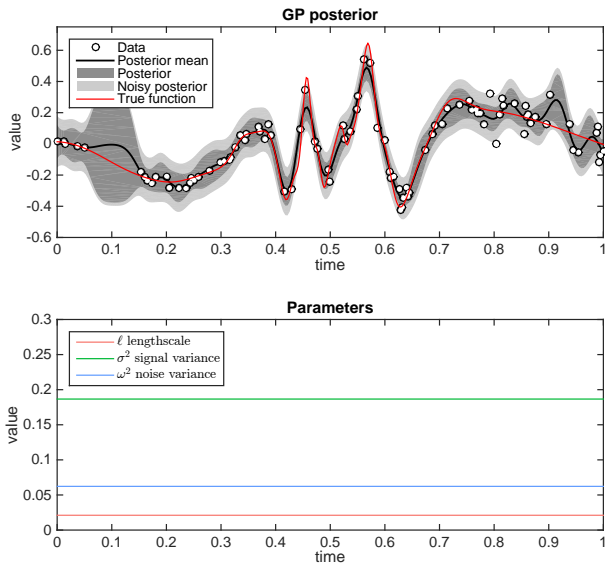
- Stationary kernels are function of vector distance  $x - x'$
- For instance if input variable is 'age' in years, then a stationary kernel has property  $K(1, 2) = K(80, 81)$
- Strange to assume that 1 and 2 year olds are **as** similar to each other as 80 and 81 year olds
- Non-stationary kernel** is not translation invariant, i.e. we can have  $K(1, 2) \neq K(80, 81)$
- Simplest non-stationary kernel is the dot product,  $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$  since
  - $\mathbf{x} = [1, 1]^T$ ,  $\mathbf{x}' = [2, 2]$ ,  $K(\mathbf{x}, \mathbf{x}') = 1 \cdot 2 + 1 \cdot 2 = 4$
  - $\mathbf{x} = [10, 10]^T$ ,  $\mathbf{x}' = [11, 11]$ ,  $K(\mathbf{x}, \mathbf{x}') = 10 \cdot 11 + 10 \cdot 11 = 120$

## Problem with stationary functions



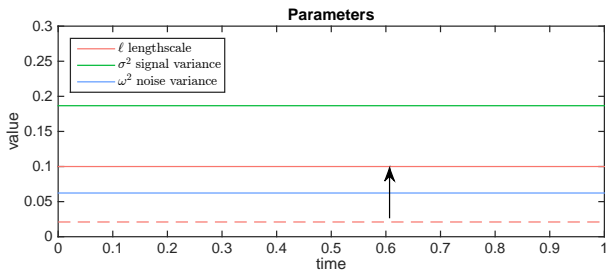
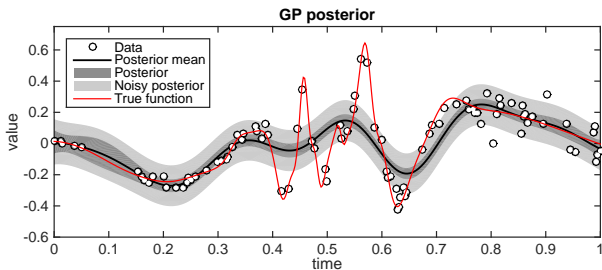
- Simple dataset

## Problem with stationary functions



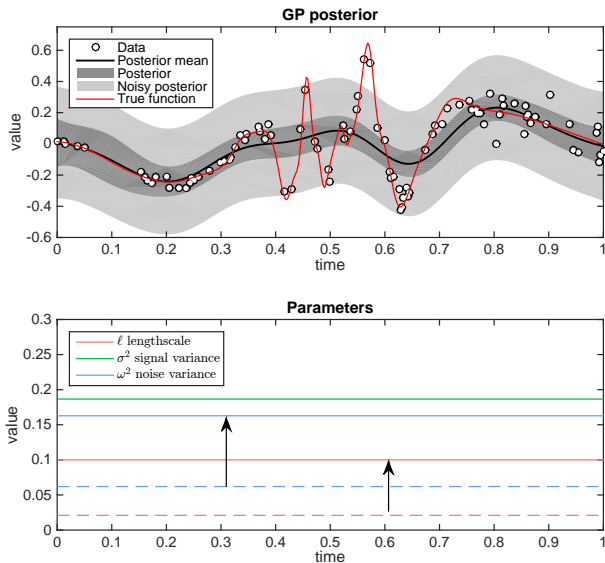
- Optimal Gaussian process fit
- Bad fit in the beginning

## Problem with stationary functions



- Let's **increase lengthscale** to get smoother model
- Initial fit fixed, now ill fit in the middle

## Problem with stationary functions



- Let's **increase noise level** to to match data
- $\Rightarrow$  We need **input-dependent** parameters

## Non-stationary Gaussian process<sup>3</sup>

- The Gaussian kernel has a fixed, global lengthscale

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right) \quad (26)$$

- Equally smooth functions everywhere
- The **non-stationary** Gaussian kernel ('Gibbs kernel') admits a lengthscale function  $\ell(x)$

$$K(x, x') = \underbrace{\sqrt{\frac{2\ell(x)\ell(x')}{\ell(x)^2 + \ell(x')^2}}}_{\text{normalizer}} \exp\left(-\frac{(x - x')^2}{\ell(x)^2 + \ell(x')^2}\right) \quad (27)$$

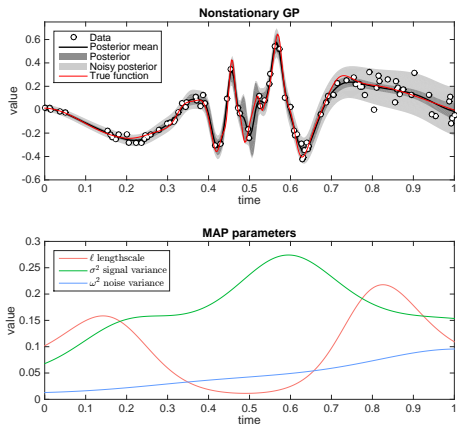
- The multivariate Gibbs kernel, where  $\Sigma_i := \Sigma(\mathbf{x}_i) \in \mathbb{R}^{D \times D}$

$$K(\mathbf{x}_i, \mathbf{x}_j) = |\Sigma_i|^{1/4} |\Sigma_j|^{1/4} |(\Sigma_i + \Sigma_j)/2|^{-1/2} \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^T ((\Sigma_i + \Sigma_j)/2)^{-1} (\mathbf{x}_i - \mathbf{x}_j)\right) \quad (28)$$

---

<sup>3</sup>Paciorek, Schervish (NIPS 2004): Nonstationary Covariance Functions for Gaussian Process Regression

## Non-stationary solution<sup>4</sup>



- Function process

$$y(x) = f(x) + \varepsilon(x) \quad (29)$$

$$f(x) \sim \mathcal{GP}(0, \sigma(x)\sigma(x')K_{\ell(\cdot)}(x, x')) \quad (30)$$

$$\varepsilon(x) \sim \mathcal{N}(0, \omega(x)^2) \quad (31)$$

- Parameter processes

$$\ell(x) \sim \mathcal{GP}(\mu_{\ell}, K_{\ell}(x, x')) \quad (32)$$

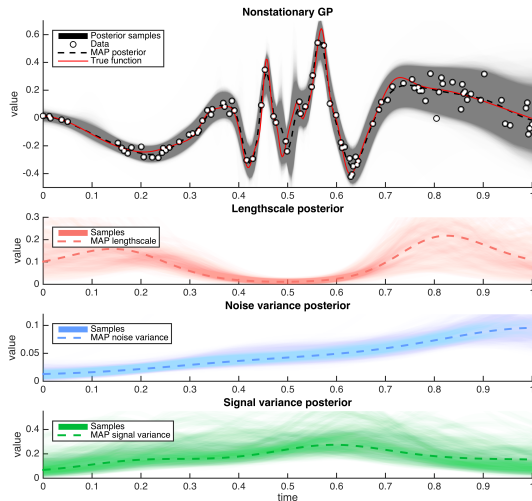
$$\sigma(x) \sim \mathcal{GP}(\mu_{\sigma}, K_{\sigma}(x, x')) \quad (33)$$

$$\omega(x) \sim \mathcal{GP}(\mu_{\omega}, K_{\omega}(x, x')) \quad (34)$$

- Kernel

$$K(x, x') = \sqrt{\frac{2\ell(x)\ell(x')}{\ell(x)^2 + \ell(x')^2}} \exp\left(-\frac{(x-x')^2}{\ell(x)^2 + \ell(x')^2}\right) \quad (35)$$

- Explicit function representation through **smoothness**, **scale** and **noise** functions



- Sample exact posterior with HMC<sup>5</sup>

$$p(\mathbf{f}, \ell, \sigma, \omega; \mathbf{y})$$

<sup>5</sup>Heinonen et al. Non-stationary Gaussian process regression with Hamiltonian Monte Carlo. AISTATS 2016



# Generalised Spectral Mixture (GSM) kernel<sup>67</sup>

- Non-stationary spectral kernel:

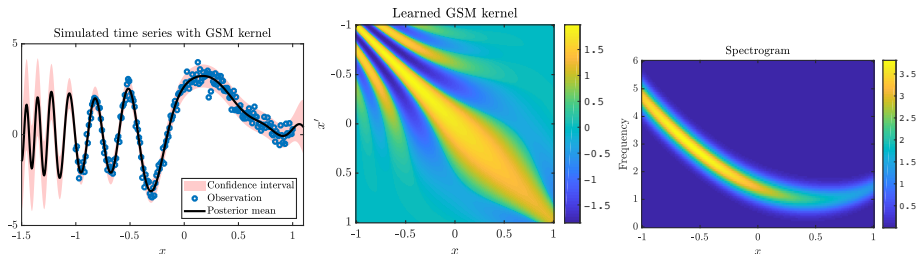
$$K_{\mathbf{w}, \mu, \sigma}(x, x') \propto \sum_{i=1}^Q \underbrace{w_i(x)w_i(x')}_{\text{Exponential kernel}} \underbrace{\exp\left(-\frac{(x-x')^2}{\ell_i(x)^2 + \ell_i(x')^2}\right) \cos(2\pi(\mu_i(x)x - \mu_i(x')x'))}_{\text{periodic}}$$

with

$$\log w_i(x) \sim \mathcal{GP}(0, K_w) \quad (36)$$

$$\log \mu_i(x) \sim \mathcal{GP}(0, K_\mu) \quad (37)$$

$$\log \ell_i(x) \sim \mathcal{GP}(0, K_\sigma) \quad (38)$$



<sup>6</sup>Remes, Heinonen, Kaski (NIPS 2017): Non-stationary spectral kernels

<sup>7</sup>Shen, Heinonen, Kaski (AISTATS 2019): Harmonizable mixture kernels with variational Fourier features

## Unified theory on spectral kernels

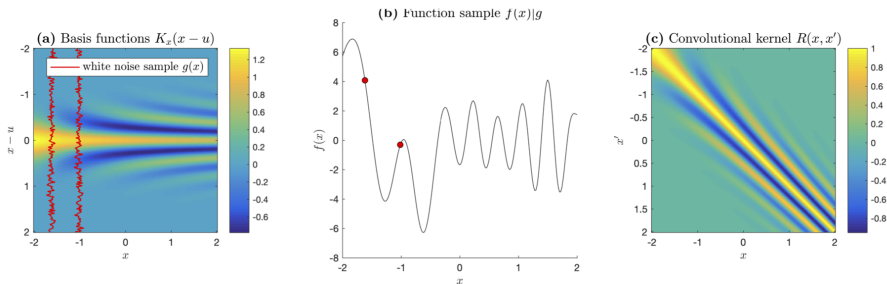
- A Gaussian process can be represented as a convolution over  $\mathbf{x}, \mathbf{u} \in \mathcal{X}$ ,

$$f(\mathbf{x}_i) = \int K_{\mathbf{x}_i}(\mathbf{u})g(\mathbf{u})d\mathbf{u} \quad (39)$$

- Feature map  $K_{\mathbf{x}_i}$
- White noise process  $g(\mathbf{u}) \sim \mathcal{GP}(0, \delta_{\mathbf{x}=\mathbf{x}'})$
- The kernel becomes<sup>8</sup>

$$C(\mathbf{x}_i, \mathbf{x}_j) = \int K_{\mathbf{x}_i}(\mathbf{u})\overline{K_{\mathbf{x}_j}(\mathbf{u})}d\mathbf{u} \quad (40)$$

where  $\overline{K}$  is complex conjugate



<sup>8</sup>Shen, Heinonen, Kaski (AISTATS 2020): Learning spectrograms with convolutional spectral kernels

## Convolutional kernel family<sup>9</sup>

- Gaussian kernel

$$K_{\mathbf{x}_i}(\mathbf{u}) \propto \mathcal{N}(\mathbf{u}|\mathbf{x}_i, \Sigma) \quad (41)$$

- A non-stationary Gaussian kernel

$$K_{\mathbf{x}_i}(\mathbf{u}) \propto \mathcal{N}(\mathbf{u}|\mathbf{x}_i, \Sigma(\mathbf{x}_i)) \quad (42)$$

- Spectral mixture kernel

$$K_{\mathbf{x}_i}(\mathbf{u}) \propto \mathcal{N}(\mathbf{u}|\mathbf{x}_i + i\boldsymbol{\mu}, \Sigma) \quad (43)$$

- Non-stationary spectral mixture kernel

$$K_{\mathbf{x}_i}(\mathbf{u}) \propto \mathcal{N}(\mathbf{u}|\mathbf{x}_i + i\boldsymbol{\mu}(\mathbf{x}_i), \Sigma(\mathbf{x}_i)) \quad (44)$$

- Input-dependent frequencies  $\boldsymbol{\mu}_i$
- Input-dependent Gaussian covariance  $\Sigma_i$

---

<sup>9</sup>Shen, Heinonen, Kaski (AISTATS 2020): Learning spectrograms with convolutional spectral kernels

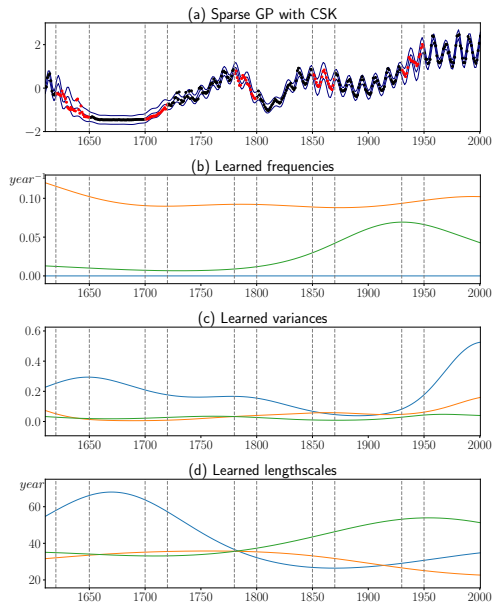
# Convolutional spectral kernel

- We parameterise frequencies with smooth GP's and learn point estimates,

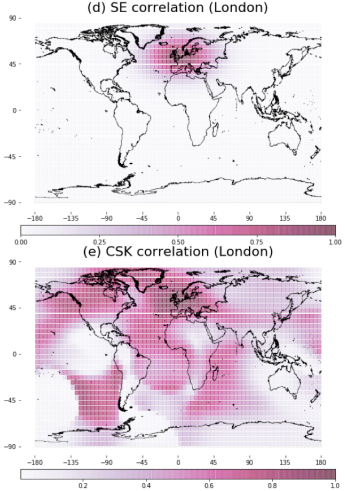
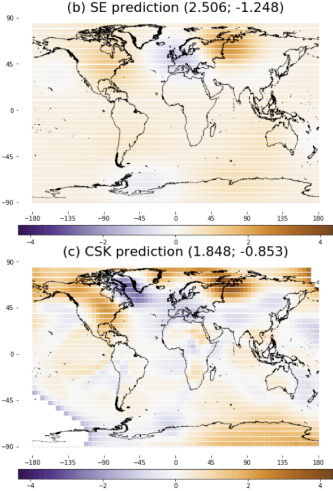
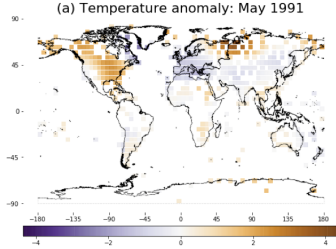
$$w^q(\cdot) \sim \mathcal{GP}(c_w, k_{SE}(\cdot, \cdot))$$

$$\text{logit} \boldsymbol{\mu}^q(\cdot) \sim \mathcal{GP}(c_\mu, k_{SE}(\cdot, \cdot))$$

$$(\boldsymbol{\Lambda}^q)^{1/2}(\cdot) \sim \mathcal{GP}(c_\lambda, k_{SE}(\cdot, \cdot))$$



# Spatial interpolation with non-stationary spectral kernel



- The non-stationary spectral kernel family is extremely flexible

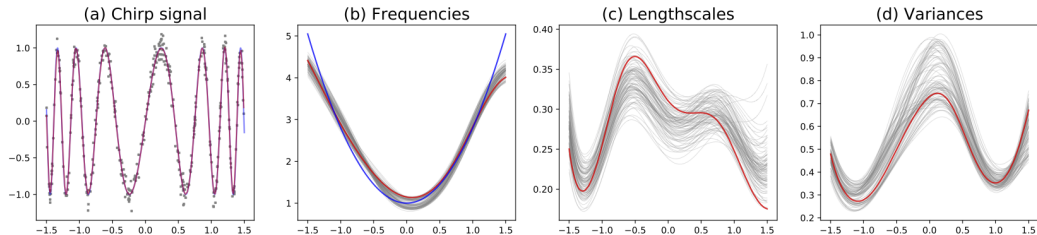
## How to learn an overly flexible kernel?

- Marginal log likelihood matches likelihood  $p(\mathbf{y}|\mathbf{f})$  and prior  $p(\mathbf{f}|\theta)$

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (45)$$

$$= \log \mathbb{E}_{p(\mathbf{f})}p(\mathbf{y}|\mathbf{f}) \quad (46)$$

- Risk overfitting if  $p(\mathbf{f}) \approx p(\mathbf{y}|\mathbf{f})$  is possible
  - Consider function space that only contains training-data like functions
- Remedies
  - MCMC sampling of  $p(\mathbf{y}|\theta)$
  - Hyperpriors  $p(\theta)$



## Summary

- The kernel choice defines how well the GP performs
- Gaussian kernel is a convenient 'default' kernel that can **interpolate** well
  - Advantage: simple, efficient, easy-to-learn, universal
  - Disadvantage: cannot fit periodic, "long-range" or non-stationary signals
- Spectral kernels can **extrapolate** repeating patterns
  - Advantage: can learn arbitrary periodic or non-periodic **stationary** patterns
  - Disadvantage: slow, possibility to overfit
- Non-stationary spectral kernels can learn **adaptive** interpolations
  - Advantage: can learn evolving frequencies
  - Disadvantage: slow, more possibilities to overfit

Thanks!





## Convolutional spectral kernel

- Convolutional spectral kernel (CSK) with Gaussian/periodic feature map

$$K_{\mathbf{x}_i}(\mathbf{u}) = \sum_{q=1}^Q w_i^q \exp(-2\pi^2 S_i^q + 2\pi i \theta_i^q) \sim \sum_{q=1}^Q \mathcal{N}(i\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad (47)$$

where  $S_i^q = (\mathbf{x}_i - \mathbf{u})^T \boldsymbol{\Lambda}_i^q (\mathbf{x}_i - \mathbf{u})$  and  $\theta_i^q = \langle \boldsymbol{\mu}_i^q, \mathbf{x}_i - \mathbf{u} \rangle$

- The kernel can be solved with  $\boldsymbol{\Sigma}_i = \boldsymbol{\Lambda}_i^{-1}$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \int_{\mathbb{R}^D} K_{\mathbf{x}_i}(\mathbf{u}) K_{\mathbf{x}_j}(\mathbf{u}) d\mathbf{u} \quad (48)$$

$$= \sum_{q,p=1}^Q \frac{w_i^p \bar{w}_j^q}{(2\pi)^{D/2} |\boldsymbol{\Lambda}_i^p + \boldsymbol{\Lambda}_j^q|^{1/2}} \exp(-\pi^2 S_{ij}^{pq} + 2\pi i \theta_{ij}^{pq} - R_{ij}^{pq}) \quad (49)$$

where

$$R_{ij}^{pq} = (\boldsymbol{\mu}_i^p - \boldsymbol{\mu}_j^q)^T ((\boldsymbol{\Lambda}_i^p + \boldsymbol{\Lambda}_j^q)/2)^{-1} (\boldsymbol{\mu}_i^p - \boldsymbol{\mu}_j^q) \quad (50)$$

$$S_{ij}^{pq} = (\mathbf{x}_i - \mathbf{x}_j)^T ((\boldsymbol{\Sigma}_i^p + \boldsymbol{\Sigma}_j^q)/2)^{-1} (\mathbf{x}_i - \mathbf{x}_j) \quad (51)$$

$$\theta_{ij}^{pq} = \langle \boldsymbol{\Lambda}_i^p (\boldsymbol{\Lambda}_i^p + \boldsymbol{\Lambda}_j^q)^{-1} \boldsymbol{\mu}_i^p + \boldsymbol{\Lambda}_j^q (\boldsymbol{\Lambda}_i^p + \boldsymbol{\Lambda}_j^q)^{-1} \boldsymbol{\mu}_j^q, \mathbf{x}_i - \mathbf{x}_j \rangle \quad (52)$$