

CS-E4075 Special course on Gaussian processes: Session #6

Markus Heinonen

Aalto University

`markus.o.heinonen@aalto.fi`

Thursday 28.1.2021

- 1 Part 1: Sneak peek to kernel theory
- 2 Part 2: Recap
- 3 Part 3: spectral kernels
- 4 Part 4: Non-stationary and heteroscedastic GPs

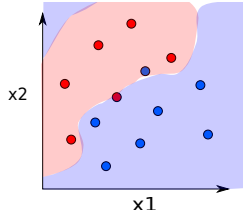
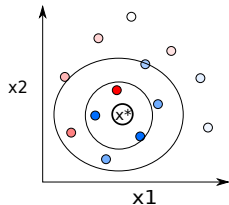
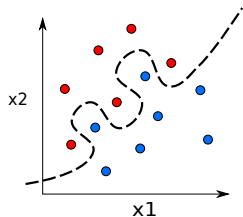
- Kernel ridge regression

$$f(\mathbf{x}) = \sum_{i=1}^N \underbrace{\alpha_i}_{\text{weight}} \underbrace{K(\mathbf{x}, \mathbf{x}_i)}_{\text{similarity}}$$

$$\boldsymbol{\alpha} = (K_{XX} + \underbrace{\lambda I}_{\text{regulariser}})^{-1} \mathbf{y} \in \mathbb{R}^N$$

- GP posterior mean coincides with KRR
- CS-E4830 Kernel Methods in Machine Learning
- Gaussian kernel (similarity)

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



- Let's study quadratic kernel over 2D inputs $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$

$$k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^2 \quad (4)$$

$$= (a_1 b_1 + a_2 b_2)^2 \quad (5)$$

$$= a_1^2 b_1^2 + 2a_1 a_2 b_1 b_2 + a_2^2 b_2^2 \quad (6)$$

$$= \underbrace{(a_1^2, \sqrt{2}a_1 a_2, a_2^2)}_{\phi(\mathbf{a})^T} \underbrace{(b_1^2, \sqrt{2}b_1 b_2, b_2^2)^T}_{\phi(\mathbf{b})}, \quad (7)$$

where $\phi(\mathbf{a}) = (a_1^2, \sqrt{2}a_1 a_2, a_2^2) \in \mathbb{R}^3$

- Linear in \mathbb{R}^3
- Non-linear in \mathbb{R}^2
- Theory: For PSD kernel there **always** exists some Hilbertian feature space \mathcal{H} such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle. \quad (8)$$

- It also means that

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i) \quad (9)$$

$$= \mathbf{w}^T \phi(\mathbf{x}), \quad \mathbf{w} \in \mathcal{H} \quad (10)$$

- Let's study quadratic kernel over 2D inputs $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$

$$k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^2 \quad (4)$$

$$= (a_1 b_1 + a_2 b_2)^2 \quad (5)$$

$$= a_1^2 b_1^2 + 2a_1 a_2 b_1 b_2 + a_2^2 b_2^2 \quad (6)$$

$$= \underbrace{(a_1^2, \sqrt{2}a_1 a_2, a_2^2)}_{\phi(\mathbf{a})^T} \underbrace{(b_1^2, \sqrt{2}b_1 b_2, b_2^2)^T}_{\phi(\mathbf{b})}, \quad (7)$$

where $\phi(\mathbf{a}) = (a_1^2, \sqrt{2}a_1 a_2, a_2^2) \in \mathbb{R}^3$

- Linear in \mathbb{R}^3
 - Non-linear in \mathbb{R}^2
- Theory: For PSD kernel there **always** exists some Hilbertian feature space \mathcal{H} such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle. \quad (8)$$

- It also means that

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i) \quad (9)$$

$$= \mathbf{w}^T \phi(\mathbf{x}), \quad \mathbf{w} \in \mathcal{H} \quad (10)$$

- Let's study quadratic kernel over 2D inputs $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$

$$k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^2 \quad (4)$$

$$= (a_1 b_1 + a_2 b_2)^2 \quad (5)$$

$$= a_1^2 b_1^2 + 2a_1 a_2 b_1 b_2 + a_2^2 b_2^2 \quad (6)$$

$$= \underbrace{(a_1^2, \sqrt{2}a_1 a_2, a_2^2)}_{\phi(\mathbf{a})^T} \underbrace{(b_1^2, \sqrt{2}b_1 b_2, b_2^2)^T}_{\phi(\mathbf{b})}, \quad (7)$$

where $\phi(\mathbf{a}) = (a_1^2, \sqrt{2}a_1 a_2, a_2^2) \in \mathbb{R}^3$

- Linear in \mathbb{R}^3
- Non-linear in \mathbb{R}^2
- Theory: For PSD kernel there **always** exists some Hilbertian feature space \mathcal{H} such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle. \quad (8)$$

- It also means that

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i) \quad (9)$$

$$= \mathbf{w}^T \phi(\mathbf{x}), \quad \mathbf{w} \in \mathcal{H} \quad (10)$$

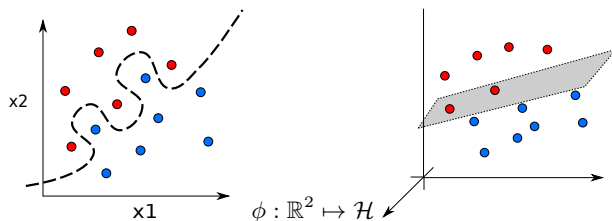
- Basis expansion ('Reproducing kernel Hilbert space, RKHS, Rasmussen 6.1.)

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \quad (11)$$

$$f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle, \quad \phi(\mathbf{x}) \in \mathcal{H} \quad (12)$$

- Gaussian kernel considers infinite number of monomials x^i

$$\phi_{gauss}(x) = e^{-x^2/2\ell^2} \left[1, \frac{1}{\sqrt{1!\ell^2}}x, \frac{1}{\sqrt{2!\ell^4}}x^2, \dots \right] \quad (13)$$



How to learn a kernel?

- Marginal likelihood: Choose a prior with maximum **amount** of functions that match the data

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$

$$\log p(\mathbf{y}|\theta) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (14)$$

$$= -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi \quad (15)$$

- ELBO: More generally variational approximation with M inducing points \mathbf{u}

$$\log p(\mathbf{y}|\theta) \geq \sum_{i=1}^N \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - KL[q(\mathbf{u})||p(\mathbf{u})] \quad (16)$$

- Relatively robust against overfitting
 - Determinant captures the volume of the data cloud in the kernel feature space
 - Finds a simple basis for the data
 - Overfitting still possible, if $p(\mathbf{f})$ contains only high-fitness solutions
- Powerful formalism to learn kernels
 - No need for model selection cross-validation
 - We can differentiate $\log p(\mathbf{y}|\theta)$ and apply gradient optimisation for parameters θ

How to learn a kernel?

- Marginal likelihood: Choose a prior with maximum **amount** of functions that match the data

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$

$$\log p(\mathbf{y}|\theta) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (14)$$

$$= -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi \quad (15)$$

- ELBO: More generally variational approximation with M inducing points \mathbf{u}

$$\log p(\mathbf{y}|\theta) \geq \sum_{i=1}^N \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - KL[q(\mathbf{u})||p(\mathbf{u})] \quad (16)$$

- Relatively robust against overfitting
 - Determinant captures the volume of the data cloud in the kernel feature space
 - Finds a simple basis for the data
 - Overfitting still possible, if $p(\mathbf{f})$ contains only high-fitness solutions
- Powerful formalism to learn kernels
 - No need for model selection cross-validation
 - We can differentiate $\log p(\mathbf{y}|\theta)$ and apply gradient optimisation for parameters θ

How to learn a kernel?

- Marginal likelihood: Choose a prior with maximum **amount** of functions that match the data

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$

$$\log p(\mathbf{y}|\theta) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (14)$$

$$= -\frac{1}{2} \underbrace{\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit}} - \frac{1}{2} \underbrace{\log |K_\theta + \sigma^2 I|}_{\text{model complexity}} - \frac{N}{2} \log 2\pi \quad (15)$$

- ELBO: More generally variational approximation with M inducing points \mathbf{u}

$$\log p(\mathbf{y}|\theta) \geq \sum_{i=1}^N \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - KL[q(\mathbf{u})||p(\mathbf{u})] \quad (16)$$

- Relatively robust against overfitting
 - Determinant captures the volume of the data cloud in the kernel feature space
 - Finds a simple basis for the data
 - Overfitting still possible, if $p(\mathbf{f})$ contains only high-fitness solutions
- Powerful formalism to learn kernels
 - No need for model selection cross-validation
 - We can differentiate $\log p(\mathbf{y}|\theta)$ and apply gradient optimisation for parameters θ

Recap (regression setting)

- 1 Gaussian process prior on inputs $\mathbf{x} \in \mathbb{R}^D$, output $y \in \mathbb{R}$,

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')) \quad (17)$$

$$\Leftrightarrow \quad (18)$$

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{m}, K_{XX}) \quad (19)$$

$$\mathbb{E}[f(\mathbf{x})] = m(\mathbf{x}) \quad (20)$$

$$\mathbf{cov}[f(\mathbf{x}), f(\mathbf{x}')] = K(\mathbf{x}, \mathbf{x}') \quad (21)$$

for inputs $X = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times D}$, functions $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^T \in \mathbb{R}^N$ and means $\mathbf{m} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))^T \in \mathbb{R}^N$,

- 2 Predictive (regression) posterior $f(\mathbf{x}) | (X, \mathbf{y}) \sim \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x})^2)$

$$\mu(\mathbf{x}) = K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} \mathbf{y} \quad (22)$$

$$\sigma(\mathbf{x})^2 = K_{\mathbf{x}\mathbf{x}} - K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} K_{X\mathbf{x}} \quad (23)$$

- 3 Optimization criteria ('- loss function') for hyperparameters θ

$$p(\mathbf{y} | \theta) = \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \theta) d\mathbf{x} = \mathcal{N}(\mathbf{y} | \mathbf{0}, K_\theta(X, X) + \sigma_n^2 I_N)$$

Recap (regression setting)

- 1 Gaussian process prior on inputs $\mathbf{x} \in \mathbb{R}^D$, output $y \in \mathbb{R}$,

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')) \quad (17)$$

$$\Leftrightarrow \quad (18)$$

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{m}, K_{XX}) \quad (19)$$

$$\mathbb{E}[f(\mathbf{x})] = m(\mathbf{x}) \quad (20)$$

$$\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = K(\mathbf{x}, \mathbf{x}') \quad (21)$$

for inputs $X = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times D}$, functions $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^T \in \mathbb{R}^N$ and means $\mathbf{m} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))^T \in \mathbb{R}^N$,

- 2 Predictive (regression) posterior $f(\mathbf{x}) | (X, \mathbf{y}) \sim \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x})^2)$

$$\mu(\mathbf{x}) = K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} \mathbf{y} \quad (22)$$

$$\sigma(\mathbf{x})^2 = K_{\mathbf{x}\mathbf{x}} - K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} K_{X\mathbf{x}} \quad (23)$$

- 3 Optimization criteria ('- loss function') for hyperparameters θ

$$p(\mathbf{y} | \theta) = \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \theta) d\mathbf{x} = \mathcal{N}(\mathbf{y} | \mathbf{0}, K_\theta(X, X) + \sigma_n^2 I_N)$$

Recap (regression setting)

- 1 Gaussian process prior on inputs $\mathbf{x} \in \mathbb{R}^D$, output $y \in \mathbb{R}$,

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')) \quad (17)$$

$$\Leftrightarrow \quad (18)$$

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{m}, K_{XX}) \quad (19)$$

$$\mathbb{E}[f(\mathbf{x})] = m(\mathbf{x}) \quad (20)$$

$$\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = K(\mathbf{x}, \mathbf{x}') \quad (21)$$

for inputs $X = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times D}$, functions $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^T \in \mathbb{R}^N$ and means $\mathbf{m} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))^T \in \mathbb{R}^N$,

- 2 Predictive (regression) posterior $f(\mathbf{x}) | (X, \mathbf{y}) \sim \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x})^2)$

$$\mu(\mathbf{x}) = K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} \mathbf{y} \quad (22)$$

$$\sigma(\mathbf{x})^2 = K_{\mathbf{x}\mathbf{x}} - K_{\mathbf{x}X} (K_{XX} + \sigma_n^2 I_N)^{-1} K_{X\mathbf{x}} \quad (23)$$

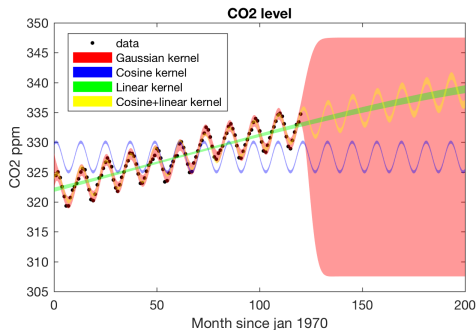
- 3 Optimization criteria ('- loss function') for hyperparameters θ

$$p(\mathbf{y} | \theta) = \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \theta) d\mathbf{x} = \mathcal{N}(\mathbf{y} | \mathbf{0}, K_\theta(X, X) + \sigma_n^2 I_N)$$

- 1 Part 1: Sneak peek to kernel theory
- 2 Part 2: Recap
- 3 Part 3: spectral kernels
- 4 Part 4: Non-stationary and heteroscedastic GPs

Which kernel to choose?

- 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 $K_{lin}(x, x') = xx' + c$
- Kernel sum $K(x, x') = K_g(x, x') + K_{lin}(x, x')$



- **Spectral kernels** can learn **arbitrary** kernel forms
 - The topic of today's lecture

Fourier transforms

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

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i x \omega} dx \quad (24)$$

where

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

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi i x \omega} d\omega \quad (25)$$

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

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + i \underbrace{\sin x}_{\text{imaginary part}} \quad (26)$$

where the complex part is often designed to cancel out (or simply ignored)

- Hence,

$$e^{\pm 2\pi i x \omega} = \cos(2\pi x \omega) \pm i \sin(2\pi x \omega) \quad (27)$$

Fourier transforms

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

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i x \omega} dx \quad (24)$$

where

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

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi i x \omega} d\omega \quad (25)$$

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

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + i \underbrace{\sin x}_{\text{imaginary part}} \quad (26)$$

where the complex part is often designed to cancel out (or simply ignored)

- Hence,

$$e^{\pm 2\pi i x \omega} = \cos(2\pi x \omega) \pm i \sin(2\pi x \omega) \quad (27)$$

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

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i x \omega} dx \quad (24)$$

where

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

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi i x \omega} d\omega \quad (25)$$

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

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + i \underbrace{\sin x}_{\text{imaginary part}} \quad (26)$$

where the complex part is often designed to cancel out (or simply ignored)

- Hence,

$$e^{\pm 2\pi i x \omega} = \cos(2\pi x \omega) \pm i \sin(2\pi x \omega) \quad (27)$$

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

$$S(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi ix\omega} dx \quad (24)$$

where

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

$$f(x) = \int_{-\infty}^{\infty} S(\omega)e^{2\pi ix\omega} d\omega \quad (25)$$

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

$$e^{ix} = \underbrace{\cos x}_{\text{real part}} + i \underbrace{\sin x}_{\text{imaginary part}} \quad (26)$$

where the complex part is often designed to cancel out (or simply ignored)

- Hence,

$$e^{\pm 2\pi ix\omega} = \cos(2\pi x\omega) \pm i \sin(2\pi x\omega) \quad (27)$$

- Let's apply Fourier to the function $K(\tau) \equiv K(x - x') = K(x, x')$, where $\tau = x - x'$

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

$$K(x - x') \equiv K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

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

- 1 All stationary kernels have *spectral density* $S(\omega)$ where ω is a frequency
 - If someone gives you a kernel $K(\tau)$, we can solve what frequencies it considers by solving the (FT)
 - Studying known kernel's frequency representations usually 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 (!)

- Let's apply Fourier to the function $K(\tau) \equiv K(x - x') = K(x, x')$, where $\tau = x - x'$

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

$$K(x - x') \equiv K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

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

- All stationary kernels have **spectral density** $S(\omega)$ where ω is a frequency
 - If someone gives you a kernel $K(\tau)$, we can solve what frequencies it considers by solving the (FT)
 - Studying known kernel's frequency representations usually of theoretical interest
- 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 (!)

- Let's apply Fourier to the function $K(\tau) \equiv K(x - x') = K(x, x')$, where $\tau = x - x'$

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

$$K(x - x') \equiv K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

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

- All stationary kernels have **spectral density** $S(\omega)$ where ω is a frequency
 - If someone gives you a kernel $K(\tau)$, we can solve what frequencies it considers by solving the (FT)
 - Studying known kernel's frequency representations usually of theoretical interest
- 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 (!)

- Let's apply Fourier to the function $K(\tau) \equiv K(x - x') = K(x, x')$, where $\tau = x - x'$

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

$$K(x - x') \equiv K(\tau) = \int_{-\infty}^{\infty} S(\omega) e^{2\pi i \omega^T \tau} d\omega \quad (\text{Inverse Fourier Transform})$$

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

- All stationary kernels have **spectral density** $S(\omega)$ where ω is a frequency
 - If someone gives you a kernel $K(\tau)$, we can solve what frequencies it considers by solving the (FT)
 - Studying known kernel's frequency representations usually of theoretical interest
- 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

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

- Hence, all stationary kernels are $S(\omega)$ -weighted combinations of sinusoids $\cos(2\pi \tau \omega)$

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

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

- Hence, all stationary kernels are $S(\omega)$ -weighted combinations of sinusoids $\cos(2\pi \tau \omega)$

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

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

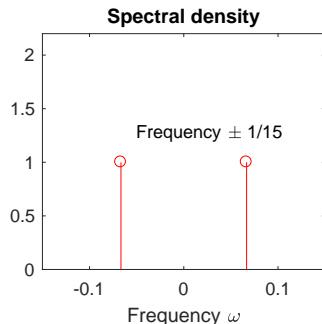
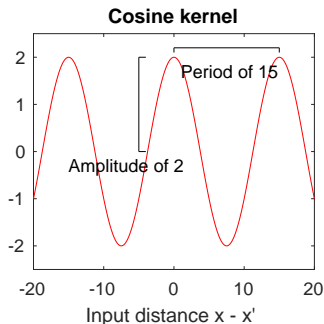
- Hence, **all** stationary kernels are $S(\omega)$ -weighted combinations of sinusoids $\cos(2\pi \tau \omega)$

Kernel sinusoid representation

- Our new general kernel definition

$$K(\tau) = \mathbb{E}_{S(\omega)} \cos(2\pi\tau\omega) \quad (28)$$

- Frequency ω is inverse of period $1/\omega$
- 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

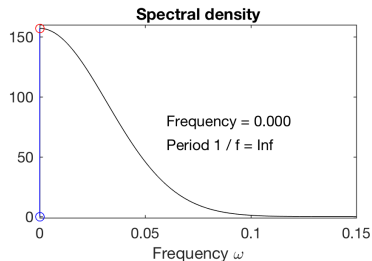
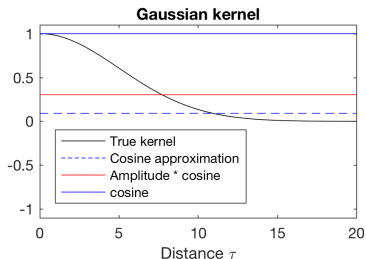
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (29)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (30)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (31)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (32)$$



Gaussian kernel sinusoids

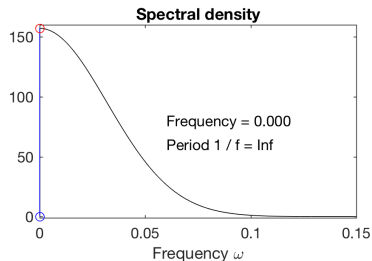
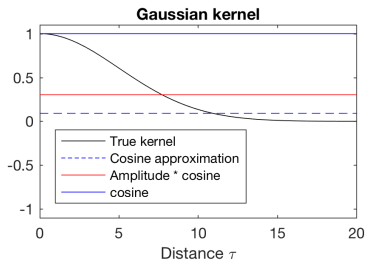
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

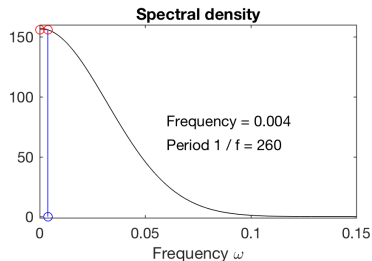
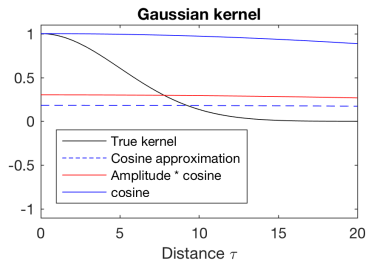
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

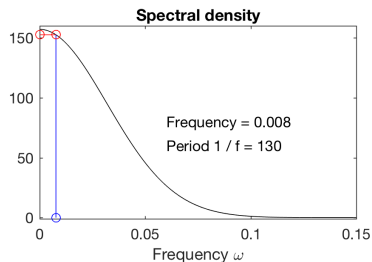
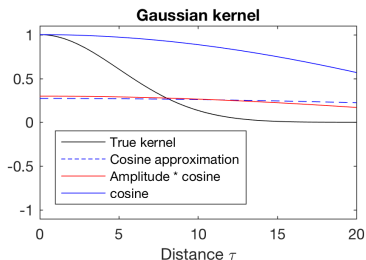
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

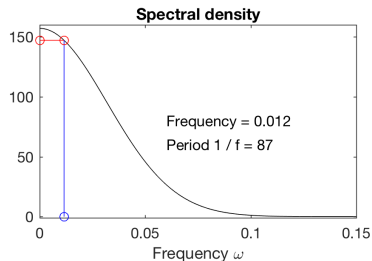
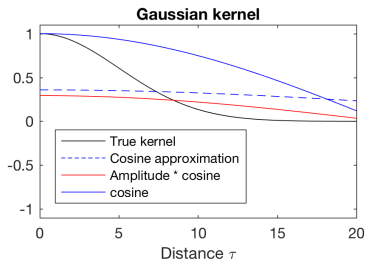
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

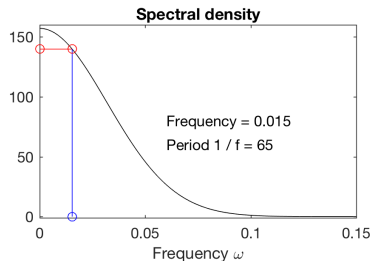
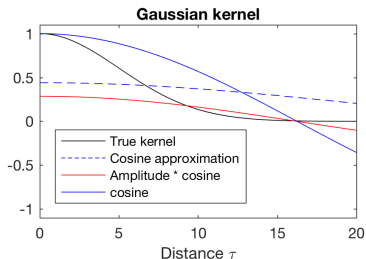
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

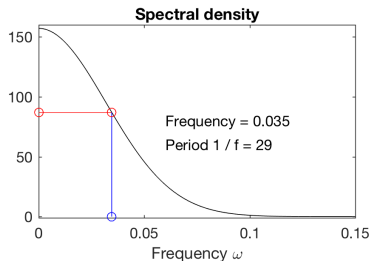
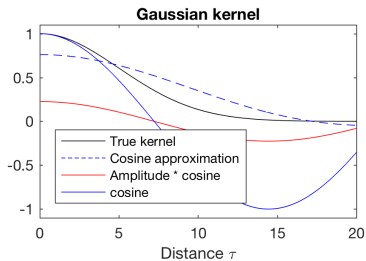
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

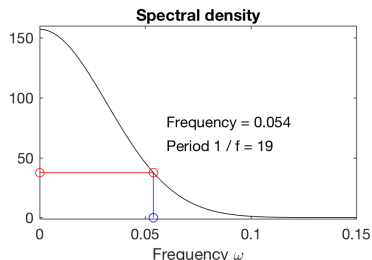
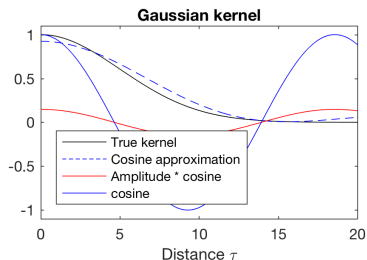
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



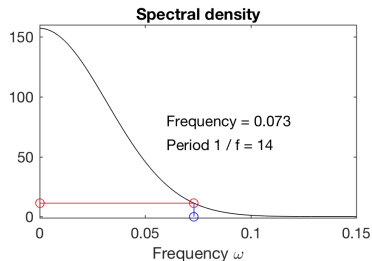
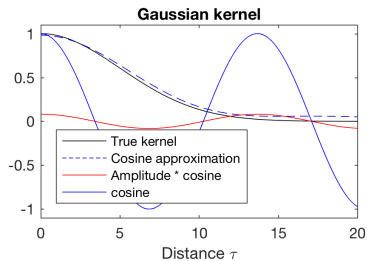
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

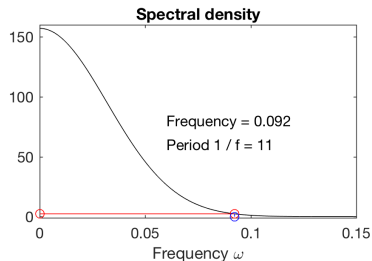
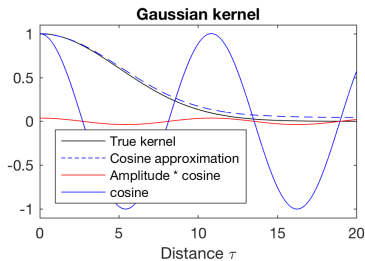
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

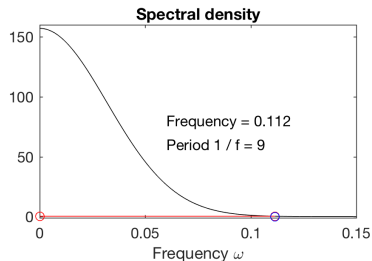
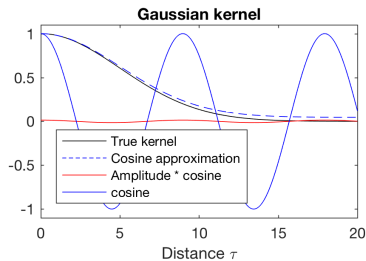
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

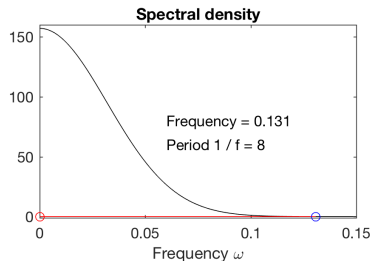
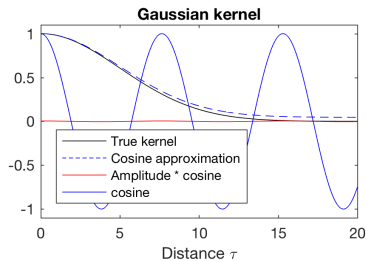
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Gaussian kernel sinusoids

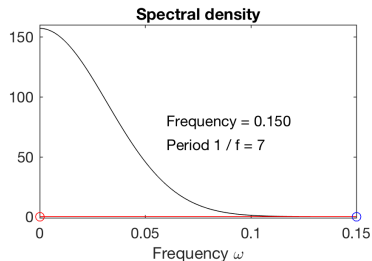
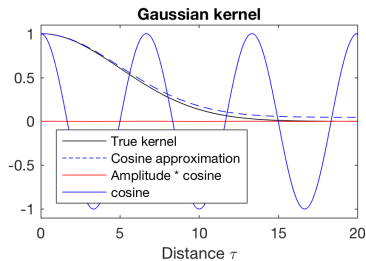
- Gaussian kernel $K_{SE}(\tau) = \exp(-\tau^2/\ell^2)$ fourier representation

$$S_{SE}(\omega) = \int_{-\infty}^{\infty} K_{SE}(\tau) e^{-2\pi i \omega^T \tau} d\tau \quad (33)$$

$$= 2\pi\ell^2 \exp(-2\pi^2\ell^2\omega^2) \quad (34)$$

$$K_{SE}(\tau) = \int_0^{\infty} \underbrace{S_{SE}(\omega)}_{\text{amplitudes}} \cdot \underbrace{\cos(2\pi\tau\omega)}_{\text{sinusoids}} d\omega \quad (35)$$

$$\approx \sum_{\omega} S_{SE}(\omega) \cdot \cos(2\pi\tau\omega) \quad (36)$$



Some spectral densities

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

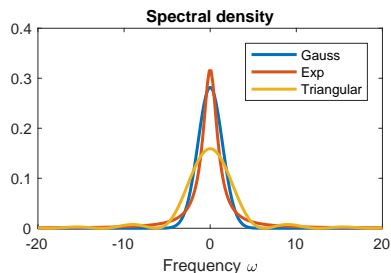
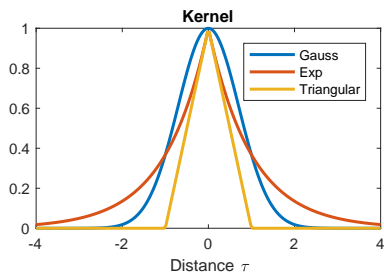
$$K_{exp}(\tau) = \exp(-|\tau|/\ell)$$

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

$$S_{gauss}(\omega) = \frac{\sqrt{\ell}}{2\sqrt{\pi}} \exp(-\ell\omega^2/4) \quad (37)$$

$$S_{exp}(\omega) = 1/(\pi/\ell + \pi\ell\omega^2) \quad (38)$$

$$S_{tri}(\omega) = (1 - \cos \omega)/(\pi\omega^2) \quad (39)$$



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

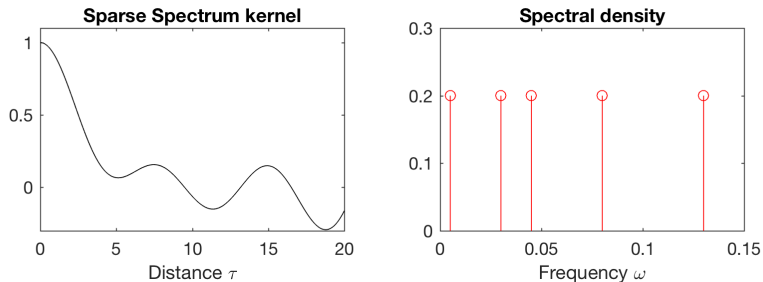
Sparse Spectrum (SS) kernel

- Define Q real frequencies $(\omega_1, \dots, \omega_Q)^T \in \mathbb{R}^Q$ with Fourier dual¹

$$S(\omega) := \frac{1}{Q} \sum_{i=1}^Q \delta(\omega - \omega_i) \quad (40)$$

$$\Rightarrow K(\tau) = \frac{1}{Q} \sum_{i=1}^Q \cos(2\pi\tau\omega_i) \quad (41)$$

- Highly structured covariance, no decay, prone to overfitting



¹Lazaro-Gredilla et al (JMLR 2010) Sparse spectrum gaussian process regression

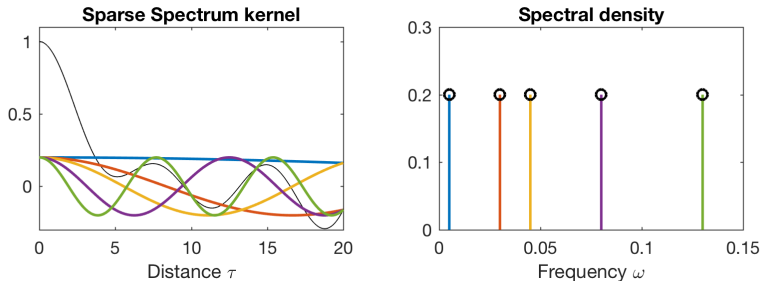
Sparse Spectrum (SS) kernel

- Define Q real frequencies $(\omega_1, \dots, \omega_Q)^T \in \mathbb{R}^Q$ with Fourier dual¹

$$S(\omega) := \frac{1}{Q} \sum_{i=1}^Q \delta(\omega - \omega_i) \quad (40)$$

$$\Rightarrow K(\tau) = \frac{1}{Q} \sum_{i=1}^Q \cos(2\pi\tau\omega_i) \quad (41)$$

- Highly structured covariance, no decay, prone to overfitting



¹Lazaro-Gredilla et al (JMLR 2010) Sparse spectrum gaussian process regression

Wilson: Spectral Mixture (SM) kernel

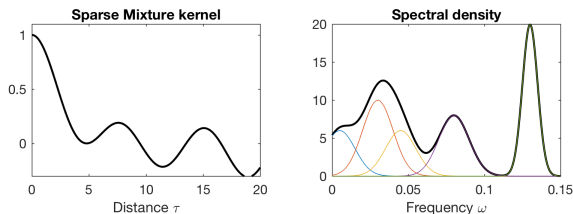
- Define mixture of Q Gaussians $\{a_i \mathcal{N}(\mu_i, \sigma_i^2)\}_{i=1}^Q$ ²

$$S(\omega) := \sum_{i=1}^Q a_i \mathcal{N}(\omega | \mu_i, \sigma_i^2) \quad (42)$$

$$\Rightarrow K(\tau) = \int_{-\infty}^{\infty} S(\omega) \cos(2\pi\tau\omega) d\omega \quad (43)$$

$$= \sum_{i=1}^Q a_i \underbrace{\exp(-2\pi^2 \sigma_i^2 \tau^2)}_{\text{smooth decay}} \underbrace{\cos(2\pi\tau\mu_i)}_{\text{periodic}} \quad (44)$$

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



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

Wilson: Spectral Mixture (SM) kernel

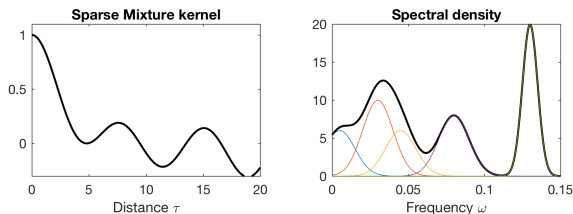
- Define mixture of Q Gaussians $\{a_i \mathcal{N}(\mu_i, \sigma_i^2)\}_{i=1}^Q$ ²

$$S(\omega) := \sum_{i=1}^Q a_i \mathcal{N}(\omega | \mu_i, \sigma_i^2) \quad (42)$$

$$\Rightarrow K(\tau) = \int_{-\infty}^{\infty} S(\omega) \cos(2\pi\tau\omega) d\omega \quad (43)$$

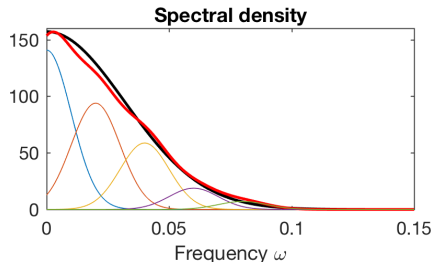
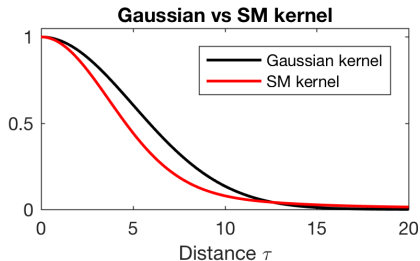
$$= \sum_{i=1}^Q a_i \underbrace{\exp(-2\pi^2 \sigma_i^2 \tau^2)}_{\text{smooth decay}} \underbrace{\cos(2\pi\tau\mu_i)}_{\text{periodic}} \quad (44)$$

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



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

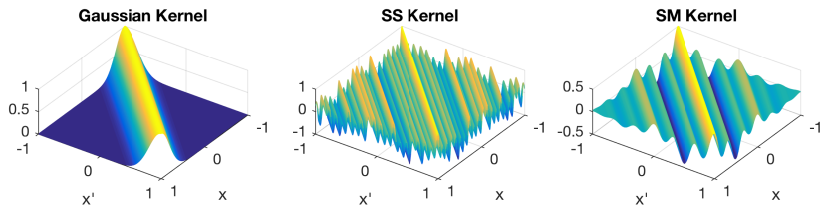
Wilson: Spectral Mixture (SM) kernel



- 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 certain a_i, μ_i, σ_i

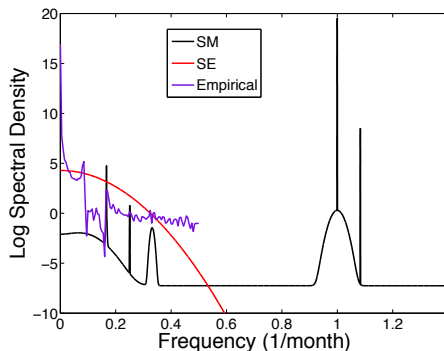
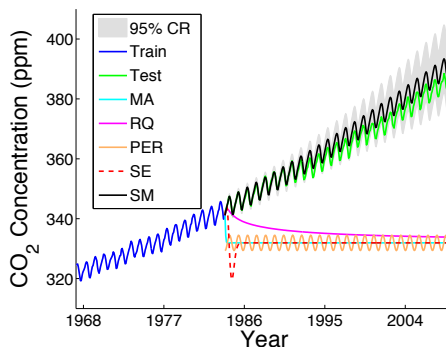


- Image from Remes, Heinonen, Kaski: Non-stationary spectral kernels, NIPS'17

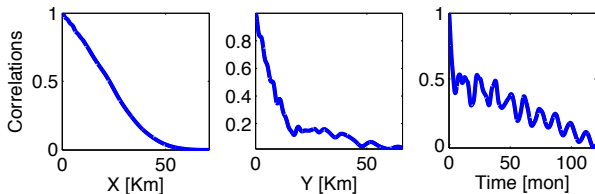
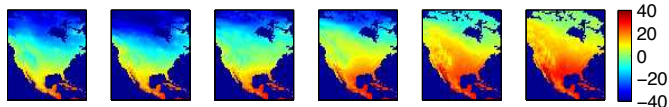
- Optimize $3Q$ hyperparameters $\theta = \{a_i, \mu_i, \sigma_i\}_{i=1}^Q$ of kernel
 $K_\theta(x - x') = \sum_{i=1}^Q a_i \exp(-2\pi^2 \sigma_i^2 \tau^2) \cos(2\pi\tau\mu_i)$ by maximizing

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

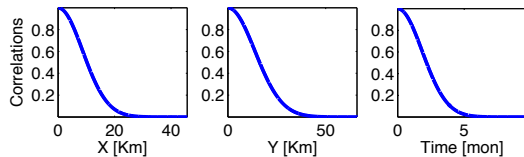
- After kernel is fixed, predictions have closed form



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

Iterative kernel learning strategies

- Deep Kernel Learning (DKL): use a neural network as a feature extractor $\text{NN} : \mathbb{R}^d \mapsto \mathbb{R}^D$

$$f(\mathbf{x}) = \mathcal{GP}(0, k(\text{NN}(\mathbf{x}), \text{NN}(\mathbf{x}'))) \quad (45)$$

- Automatic Statistician / Automatic Bayesian Covariance Discovery (ABCD) / Neural Kernel Network (NKN): Search over kernel sums and products

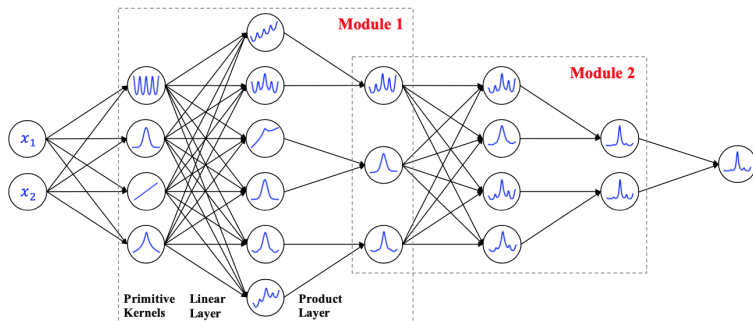


Figure 2. Neural Kernel Network: each module consists of a **Linear** layer and a **Product** layer. NKN is based on compositional rules for kernels, thus every individual unit itself represents a kernel.

- 1 Part 1: Sneak peek to kernel theory
- 2 Part 2: Recap
- 3 Part 3: spectral kernels
- 4 Part 4: Non-stationary and heteroscedastic GPs

- Standard GP regression

$$y(\mathbf{x}) = f(\mathbf{x}) + \varepsilon \quad (46)$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2) \quad (47)$$

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k_\ell(\mathbf{x}, \mathbf{x}')) \quad (48)$$

- Global noise variance
 - Global kernel function
- Heteroscedastic GPs: What if noise depends on inputs?
 - Non-stationary GPs: What if function dynamics depends on inputs?

- Standard GP regression

$$y(\mathbf{x}) = f(\mathbf{x}) + \varepsilon \quad (46)$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2) \quad (47)$$

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k_\ell(\mathbf{x}, \mathbf{x}')) \quad (48)$$

- Global noise variance
- Global kernel function
- Heteroscedastic GPs: What if noise depends on inputs?
- Non-stationary GPs: What if function dynamics depends on inputs?

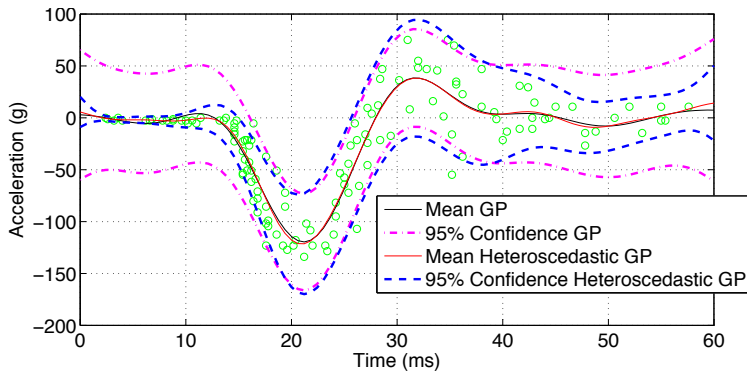


Figure 1. Silverman's (1985) motorcycle benchmark is an example for input dependent noise. It consists of a sequence of accelerometer readings through time following a simulated motor-cycle crash.

³Kersting et al (2007): Most Likely Heteroscedastic Gaussian process regression

- Standard Gaussian process assumes **additive zero-mean noise model**

$$y(\mathbf{x}) = f(\mathbf{x}) + \varepsilon(\mathbf{x}) \quad (49)$$

$$\varepsilon(\mathbf{x}) \sim \mathcal{N}(0, \sigma_n^2) \quad (50)$$

where all noises are zero mean with constant variance σ_n^2

- Heteroscedastic model assumes **input-dependent** noise:

$$\varepsilon(\mathbf{x}) \sim \mathcal{N}(0, \sigma_n(\mathbf{x})^2)$$

- More complex (non-Gaussian) noise models are sometimes used
- The function $\sigma_n(\mathbf{x})^2$ can be another Gaussian process (!)
- Leads to a joint model

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

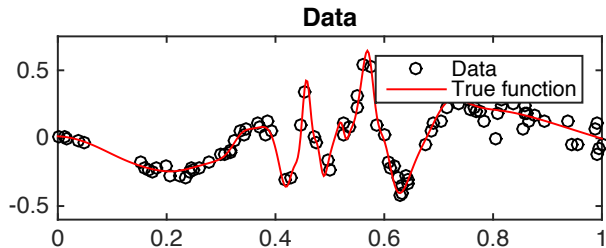
$$K(x, x') = K(x + a, x' + a) \quad (51)$$

$$K(x, x') = K(x - x') \quad (52)$$

for any a

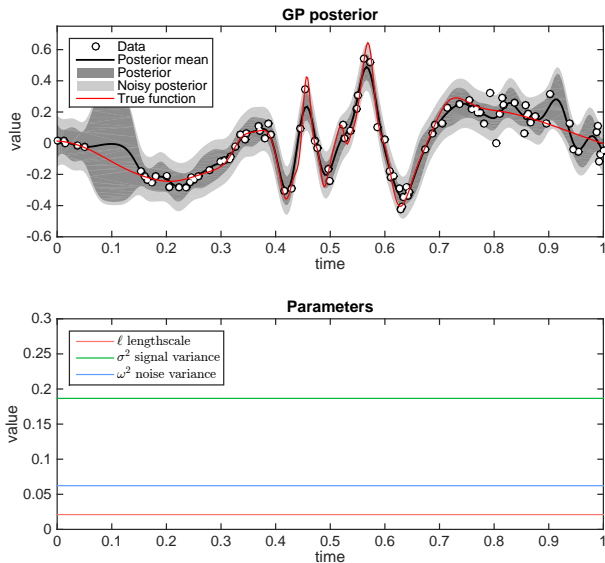
- Stationary kernels are function of vector distance $x - x'$
- For instance if input variable is 'age' in years, then a stationary kernel has property $K(1, 2) = K(80, 81)$
- Strange to assume that 1 and 2 year olds are **as** similar to each other as 80 and 81 year olds
- **Non-stationary kernel** 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



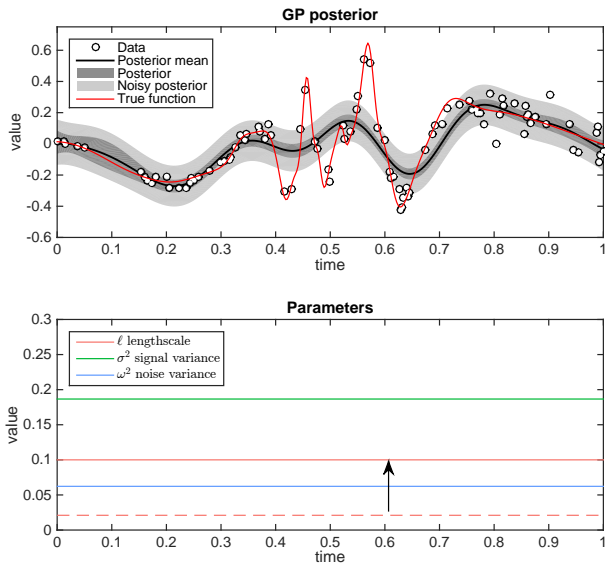
- Simple dataset

Problem with stationary functions



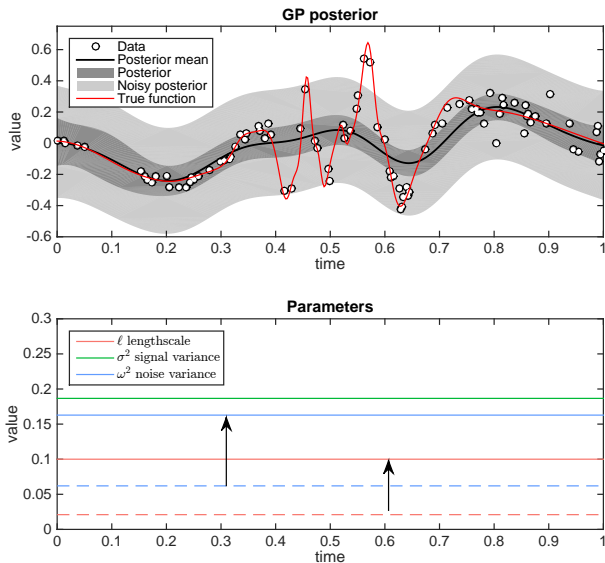
- Optimal Gaussian process fit
- Bad fit in the beginning

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

- The Gaussian kernel has a fixed, global lengthscale

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right) \quad (53)$$

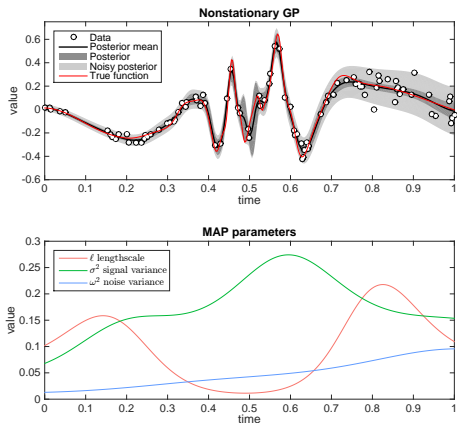
- Equally smooth functions everywhere
- The **non-stationary** Gaussian kernel ('Gibbs kernel') admits a lengthscale function $\ell(x)$

$$K(x, x') = \underbrace{\sqrt{\frac{2\ell(x)\ell(x')}{\ell(x)^2 + \ell(x')^2}}}_{\text{normalizer}} \exp\left(-\frac{(x - x')^2}{\ell(x)^2 + \ell(x')^2}\right) \quad (54)$$

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

$$K(\mathbf{x}_i, \mathbf{x}_j) = |\Sigma_i|^{1/4} |\Sigma_j|^{1/4} |(\Sigma_i + \Sigma_j)/2|^{-1/2} \exp\left(-(\mathbf{x}_i - \mathbf{x}_j)^T ((\Sigma_i + \Sigma_j)/2)^{-1} (\mathbf{x}_i - \mathbf{x}_j)\right) \quad (55)$$

Non-stationary solution⁴



- Function process

$$y(x) = f(x) + \varepsilon(x) \quad (56)$$

$$f(x) \sim \mathcal{GP}(0, \sigma(x)\sigma(x')K_{\ell(\cdot)}(x, x')) \quad (57)$$

$$\varepsilon(x) \sim \mathcal{N}(0, \omega(x)^2) \quad (58)$$

- Parameter processes

$$\ell(x) \sim \mathcal{GP}(\mu_{\ell}, K_{\ell}(x, x')) \quad (59)$$

$$\sigma(x) \sim \mathcal{GP}(\mu_{\sigma}, K_{\sigma}(x, x')) \quad (60)$$

$$\omega(x) \sim \mathcal{GP}(\mu_{\omega}, K_{\omega}(x, x')) \quad (61)$$

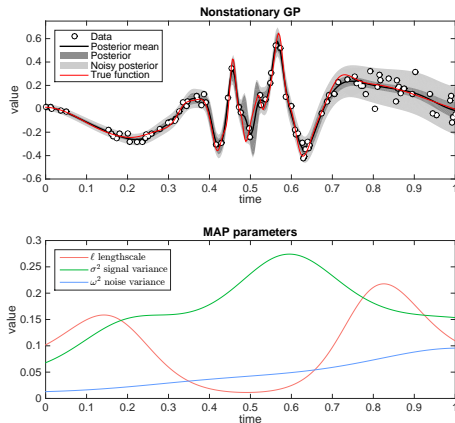
- Kernel

$$K(x, x') = \sqrt{\frac{2\ell(x)\ell(x')}{\ell(x)^2 + \ell(x')^2}} \exp\left(-\frac{(x-x')^2}{\ell(x)^2 + \ell(x')^2}\right) \quad (62)$$

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

⁴Heinonen et al. Non-stationary Gaussian process regression with Hamiltonian Monte Carlo. AISTATS 2016

Non-stationary inference

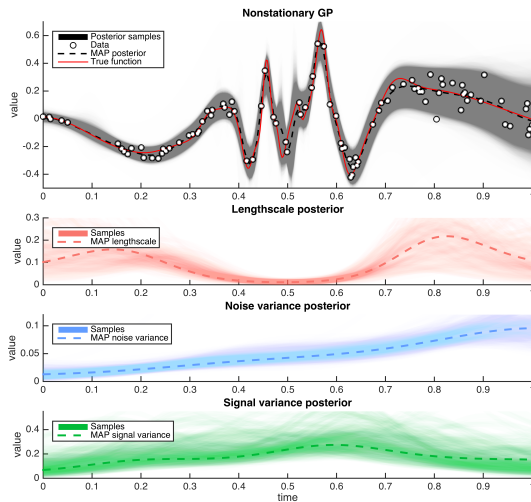


- Marginal joint likelihood

$$\mathcal{L} = p(\mathbf{y}, \ell, \omega, \sigma) = p(\mathbf{y}|\ell, \omega, \sigma)p(\ell)p(\sigma)p(\omega) \quad (63)$$

$$= \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma\sigma^T \circ K_\ell + \text{diag}(\omega))\mathcal{N}(\ell|\mu_\ell, K_\ell)\mathcal{N}(\sigma|\mu_\sigma, K_\sigma)\mathcal{N}(\omega|\mu_\omega, K_\omega) \quad (64)$$

- We optimize \mathcal{L} for MAP estimates $\hat{\ell}, \hat{\sigma}, \hat{\omega}$.
- The predictive posterior $p(\mathbf{f}|\hat{\ell}, \hat{\sigma}, \hat{\omega}, \mathbf{y})$ is of standard form, except our kernel is $\hat{\sigma}\hat{\sigma}^T \circ K_{\hat{\ell}}$



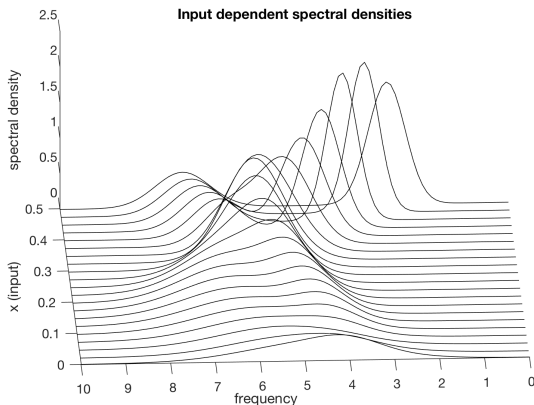
- Sample exact posterior with HMC⁵

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

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

Non-stationary spectral kernels

- We have seen how to learn arbitrary **stationary** kernels via spectral learning
- We have seen how to learn (non-stationary) Gaussian kernel with parameter functions
- What about non-stationary spectral kernels?
- Model input-dependent frequencies, or spectrograms $S(x, \omega)$
 - E.g. wavelets are time-dependent frequencies in signal processing



Generalised Spectral Mixture (GSM) kernel⁶⁷

- Non-stationary spectral kernel can be derived:

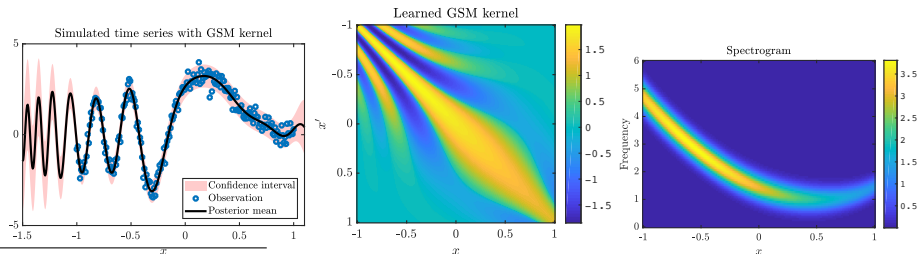
$$K_{\mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\sigma}}(x, x') \propto \sum_{i=1}^Q \underbrace{w_i(x)w_i(x')}_{\text{Exponential kernel}} \underbrace{\exp\left(-\frac{(x-x')^2}{\ell_i(x)^2 + \ell_i(x')^2}\right) \cos(2\pi(\mu_i(x)x - \mu_i(x')x'))}_{\text{periodic}}$$

with

$$\log w_i(x) \sim \mathcal{GP}(0, K_w) \quad (65)$$

$$\log \mu_i(x) \sim \mathcal{GP}(0, K_\mu) \quad (66)$$

$$\log \ell_i(x) \sim \mathcal{GP}(0, K_\sigma) \quad (67)$$



⁶Remes, Heinonen, Kaski (2017): Non-stationary spectral kernels

⁷Shen, Heinonen, Kaski (2019): Harmonizable mixture kernels with variational Fourier features

- Performance of GP has **crucial** dependency on how well the kernel matches the data
- Gaussian kernel is a convenient 'default' kernel that can **interpolate** well
 - Advantage: simple, efficient, easy-to-learn, universal
 - Disadvantage: cannot fit periodic data, stationary only
- Spectral kernels can **extrapolate** repeating patterns
 - Advantage: can learn arbitrary periodic or non-periodic **stationary** patterns
 - Disadvantage: slower to learn, high possibility to overfit
- Non-stationary Gaussian kernel can learn **adaptive** interpolations
 - Advantage: can learn smoothly changing smoothness / variance
 - Disadvantage: slower to learn, more possibilities to overfit
- Compositional kernels search for base kernel combinations

- Performance of GP has **crucial** dependency on how well the kernel matches the data
- Gaussian kernel is a convenient 'default' kernel that can **interpolate** well
 - Advantage: simple, efficient, easy-to-learn, universal
 - Disadvantage: cannot fit periodic data, stationary only
- Spectral kernels can **extrapolate** repeating patterns
 - Advantage: can learn arbitrary periodic or non-periodic **stationary** patterns
 - Disadvantage: slower to learn, high possibility to overfit
- Non-stationary Gaussian kernel can learn **adaptive** interpolations
 - Advantage: can learn smoothly changing smoothness / variance
 - Disadvantage: slower to learn, more possibilities to overfit
- Compositional kernels search for base kernel combinations

- Performance of GP has **crucial** dependency on how well the kernel matches the data
- Gaussian kernel is a convenient 'default' kernel that can **interpolate** well
 - Advantage: simple, efficient, easy-to-learn, universal
 - Disadvantage: cannot fit periodic data, stationary only
- Spectral kernels can **extrapolate** repeating patterns
 - Advantage: can learn arbitrary periodic or non-periodic **stationary** patterns
 - Disadvantage: slower to learn, high possibility to overfit
- Non-stationary Gaussian kernel can learn **adaptive** interpolations
 - Advantage: can learn smoothly changing smoothness / variance
 - Disadvantage: slower to learn, more possibilities to overfit
- Compositional kernels search for base kernel combinations