

# Probabilistic Graphical Models

## Lecture 21: Advanced Gaussian Processes

Andrew Gordon Wilson

[www.cs.cmu.edu/~andrewgw](http://www.cs.cmu.edu/~andrewgw)  
Carnegie Mellon University



April 1, 2015

# Gaussian process review

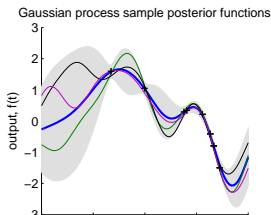
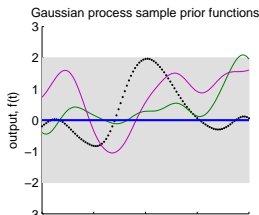
## Definition

A Gaussian process (GP) is a collection of random variables, any finite number of which have a joint Gaussian distribution.

## Nonparametric Regression Model

- Prior:  $f(x) \sim \mathcal{GP}(m(x), k(x, x'))$ , meaning  $(f(x_1), \dots, f(x_N)) \sim \mathcal{N}(\boldsymbol{\mu}, K)$ , with  $\boldsymbol{\mu}_i = m(x_i)$  and  $K_{ij} = \text{cov}(f(x_i), f(x_j)) = k(x_i, x_j)$ .

$$\overbrace{p(f(x)|\mathcal{D})}^{\text{GP posterior}} \propto \overbrace{p(\mathcal{D}|f(x))}^{\text{Likelihood}} \overbrace{p(f(x))}^{\text{GP prior}}$$



# Gaussian Process Inference

- ▶ Observed noisy data  $\mathbf{y} = (y(x_1), \dots, y(x_N))^T$  at input locations  $X$ .
- ▶ Start with the standard regression assumption:  $\mathcal{N}(y(x); f(x), \sigma^2)$ .
- ▶ Place a Gaussian process distribution over noise free functions  $f(x) \sim \mathcal{GP}(0, k_\theta)$ . The kernel  $k$  is parametrized by  $\theta$ .
- ▶ Infer  $p(\mathbf{f}_* | \mathbf{y}, X, X_*)$  for the noise free function  $f$  evaluated at test points  $X_*$ .

## Joint distribution

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K_\theta(X, X) + \sigma^2 I & K_\theta(X, X_*) \\ K_\theta(X_*, X) & K_\theta(X_*, X_*) \end{bmatrix} \right). \quad (1)$$

## Conditional predictive distribution

$$\mathbf{f}_* | X_*, X, \mathbf{y}, \theta \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \quad (2)$$

$$\bar{\mathbf{f}}_* = K_\theta(X_*, X) [K_\theta(X, X) + \sigma^2 I]^{-1} \mathbf{y}, \quad (3)$$

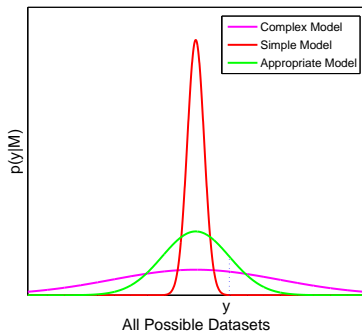
$$\text{cov}(\mathbf{f}_*) = K_\theta(X_*, X_*) - K_\theta(X_*, X) [K_\theta(X, X) + \sigma^2 I]^{-1} K_\theta(X, X_*). \quad (4)$$

# Learning and Model Selection

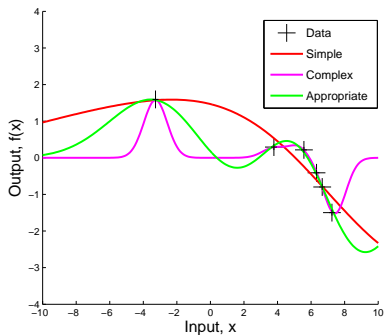
$$p(\mathcal{M}_i|\mathbf{y}) = \frac{p(\mathbf{y}|\mathcal{M}_i)p(\mathcal{M}_i)}{p(\mathbf{y})} \quad (5)$$

We can write the *evidence* of the model as

$$p(\mathbf{y}|\mathcal{M}_i) = \int p(\mathbf{y}|\mathbf{f}, \mathcal{M}_i)p(\mathbf{f})d\mathbf{f}, \quad (6)$$



(a)



(b)

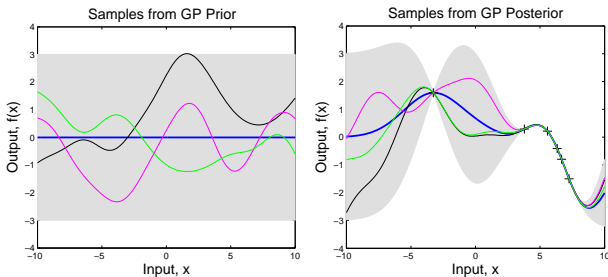
# Learning and Model Selection

- We can integrate away the entire Gaussian process  $f(x)$  to obtain the marginal likelihood, as a function of kernel hyperparameters  $\theta$  alone.

$$p(\mathbf{y}|\boldsymbol{\theta}, X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|\boldsymbol{\theta}, X)d\mathbf{f}. \quad (7)$$

$$\log p(\mathbf{y}|\boldsymbol{\theta}, X) = \underbrace{-\frac{1}{2}\mathbf{y}^T(K_{\boldsymbol{\theta}} + \sigma^2 I)^{-1}\mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2}\log |K_{\boldsymbol{\theta}} + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2}\log(2\pi). \quad (8)$$

- An extremely powerful mechanism for kernel learning.



# Inference and Learning

1. **Learning:** Optimize marginal likelihood,

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

with respect to kernel hyperparameters  $\boldsymbol{\theta}$ .

2. **Inference:** Conditioned on kernel hyperparameters  $\boldsymbol{\theta}$ , form the predictive distribution for test inputs  $X_*$ :

$$\mathbf{f}_*|X_*, X, \mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)) ,$$

$$\bar{\mathbf{f}}_* = K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}\mathbf{y} ,$$

$$\text{cov}(\mathbf{f}_*) = K_{\boldsymbol{\theta}}(X_*, X_*) - K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}K_{\boldsymbol{\theta}}(X, X_*) .$$

# Learning and Model Selection

- ▶ A fully Bayesian treatment would integrate away kernel hyperparameters  $\theta$ .

$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) = \int p(\mathbf{f}_* | X_*, X, \mathbf{y}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \quad (9)$$

- ▶ For example, we could specify a prior  $p(\boldsymbol{\theta})$ , use MCMC to take  $J$  samples from  $p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$ , and then find

$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) \approx \frac{1}{J} \sum_{i=1}^J p(\mathbf{f}_* | X_*, X, \mathbf{y}, \boldsymbol{\theta}^{(i)}), \quad \boldsymbol{\theta}^{(i)} \sim p(\boldsymbol{\theta} | \mathbf{y}). \quad (10)$$

- ▶ If we have a non-Gaussian noise model, and thus cannot integrate away  $\mathbf{f}$ , the strong dependencies between Gaussian process  $\mathbf{f}$  and hyperparameters  $\boldsymbol{\theta}$  make sampling extremely difficult. In my experience, the most effective solution is to use a deterministic approximation for the posterior  $p(\mathbf{f} | \mathbf{y})$  which enables one to work with an approximate marginal likelihood.

Let  $\tau = x - x'$ :

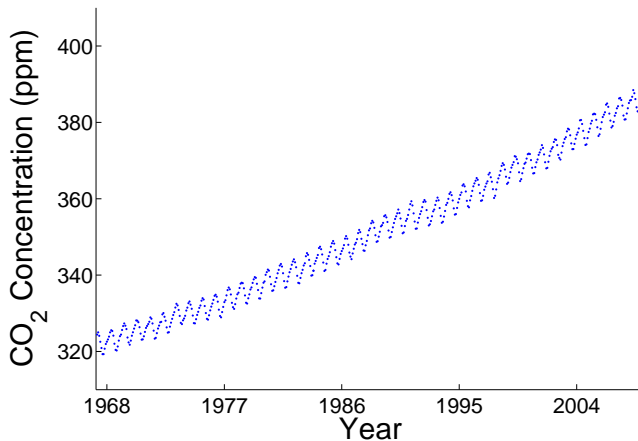
$$k_{\text{SE}}(\tau) = \exp(-0.5\tau^2/\ell^2) \quad (11)$$

$$k_{\text{MA}}(\tau) = a(1 + \frac{\sqrt{3}\tau}{\ell}) \exp(-\frac{\sqrt{3}\tau}{\ell}) \quad (12)$$

$$k_{\text{RQ}}(\tau) = (1 + \frac{\tau^2}{2\alpha\ell^2})^{-\alpha} \quad (13)$$

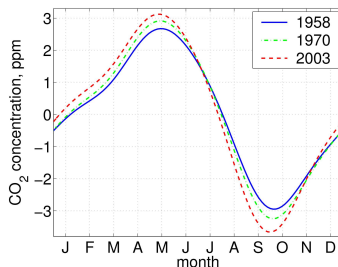
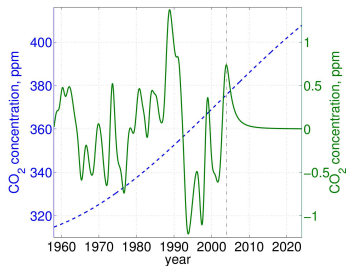
$$k_{\text{PE}}(\tau) = \exp(-2 \sin^2(\pi \tau \omega)/\ell^2) \quad (14)$$

# Worked Example: Combining Kernels, CO<sub>2</sub> Data

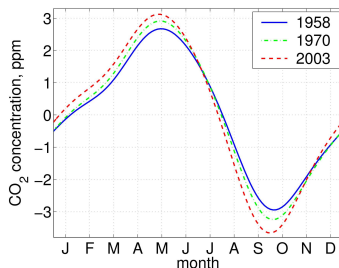
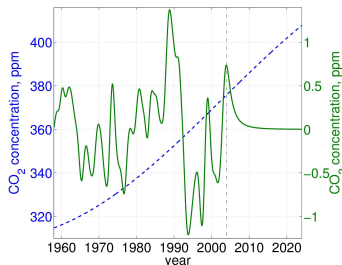


Example from Rasmussen and Williams (2006), *Gaussian Processes for Machine Learning*.

# Worked Example: Combining Kernels, CO<sub>2</sub> Data



# Worked Example: Combining Kernels, CO<sub>2</sub> Data



- ▶ Long rising trend:  $k_1(x_p, x_q) = \theta_1^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_2^2}\right)$
- ▶ Quasi-periodic seasonal changes:  $k_2(x_p, x_q) = k_{\text{RBF}}(x_p, x_q)k_{\text{PER}}(x_p, x_q) = \theta_3^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_4^2} - \frac{2 \sin^2(\pi(x_p - x_q))}{\theta_5^2}\right)$
- ▶ Multi-scale medium term irregularities:  
$$k_3(x_p, x_q) = \theta_6^2 \left(1 + \frac{(x_p - x_q)^2}{2\theta_8\theta_7^2}\right)^{-\theta_8}$$
- ▶ Correlated and i.i.d. noise:  $k_4(x_p, x_q) = \theta_9^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_{10}^2}\right) + \theta_{11}^2 \delta_{pq}$
- ▶  $k_{\text{total}}(x_p, x_q) = k_1(x_p, x_q) + k_2(x_p, x_q) + k_3(x_p, x_q) + k_4(x_p, x_q)$

# What is a kernel?

- ▶ Informally,  $k$  describes the similarities between pairs of data points. For example, far away points may be considered less similar than nearby points.  $K_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$  and so tells us the overlap between the features (basis functions)  $\phi(x_i)$  and  $\phi(x_j)$
- ▶ We have seen that all linear basis function models  $f(x) = \mathbf{w}^T \phi(x)$ , with  $p(\mathbf{w}) = \mathcal{N}(0, \Sigma_w)$  correspond to Gaussian processes with kernel  $k(x, x') = \phi(x)^T \Sigma_w \phi(x')$ .
- ▶ We have also accumulated some experience with the RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x-x'\|^2}{2\ell^2})$ .
- ▶ The kernel controls the generalisation behaviour of a kernel machine. For example, a kernel controls the support and inductive biases of a Gaussian process – which functions are a priori likely.
- ▶ A kernel is also known as covariance function or covariance kernel in the context of Gaussian processes.

# Candidate Kernel

$$k(x, x') = \begin{cases} 1 & \|x - x'\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- ▶ Symmetric
- ▶ Provides information about proximity of points
- ▶ Exercise: Is it a valid kernel?

# Candidate Kernel

$$k(x, x') = \begin{cases} 1 & \|x - x'\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Try the points  $x_1 = 1, x_2 = 2, x_3 = 3$ .

Compute the kernel matrix

$$K = \begin{bmatrix} ? & ? & ? \\ ? & ? & ? \\ ? & ? & ? \end{bmatrix} \quad (15)$$

# Candidate Kernel

$$k(x, x') = \begin{cases} 1 & \|x - x'\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Try the points  $x_1 = 1, x_2 = 2, x_3 = 3$ .

Compute the kernel matrix

$$K = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix} \quad (16)$$

The eigenvalues of  $K$  are  $(\sqrt{2} - 1)^{-1}$ , 1, and  $(1 - \sqrt{2})$ . Therefore  $K$  is not positive semidefinite.

# Representer Theorem

A decision function  $f(x)$  can be written as

$$f(x) = \langle \mathbf{w}, \phi(x) \rangle = \left\langle \sum_{i=1}^N \alpha_i \phi(x_i), \phi(x) \right\rangle = \sum_{i=1}^N \alpha_i k(x_i, x). \quad (17)$$

- ▶ Representer theorem says this function exists with finitely many coefficients  $\alpha_i$  even when  $\phi$  is infinite dimensional (an infinite number of basis functions).
- ▶ Initially viewed as a strength of kernel methods, for datasets not exceeding e.g. ten thousand points.
- ▶ Unfortunately, the number of nonzero  $\alpha_i$  often grows linearly in the size of the training set  $N$ .
- ▶ Example: In GP regression, the predictive mean is

$$\mathbb{E}[f_* | \mathbf{y}, X, x_*] = \mathbf{k}_*^T (K + \sigma^2 I)^{-1} \mathbf{y} = \sum_{i=1}^N \alpha_i k(x_i, x_*), \quad (18)$$

where  $\alpha_i = (K + \sigma^2 I)^{-1} \mathbf{y}$ .

# Making new kernels from old

Suppose  $k_1(x, x')$  and  $k_2(x, x')$  are valid. Then the following covariance functions are also valid:

$$k(x, x') = g(x)k_1(x, x')g(x') \quad (19)$$

$$k(x, x') = q(k_1(x, x')) \quad (20)$$

$$k(x, x') = \exp(k_1(x, x')) \quad (21)$$

$$k(x, x') = k_1(x, x') + k_2(x, x') \quad (22)$$

$$k(x, x') = k_1(x, x')k_2(x, x') \quad (23)$$

$$k(x, x') = k_3(\phi(x), \phi(x')) \quad (24)$$

$$k(x, x') = x^T A x' \quad (25)$$

$$k(x, x') = k_a(x_a, x'_a) + k_b(x_b, x'_b) \quad (26)$$

$$k(x, x') = k_a(x_a, x'_a)k_b(x_b, x'_b) \quad (27)$$

where  $g$  is any function,  $q$  is a polynomial with nonnegative coefficients,  $\phi(x)$  is a function from  $x$  to  $\mathbb{R}^M$ ,  $k_3$  is a valid covariance function in  $\mathbb{R}^M$ ,  $A$  is a symmetric positive definite matrix,  $x_a$  and  $x_b$  are not necessarily disjoint variables with  $x = (x_a, x_b)^T$ , and  $k_a$  and  $k_b$  are valid kernels in their respective spaces.

# Stationary Kernels

- ▶ A *stationary* kernel is invariant to translations of the input space. Equivalently,  $k = k(x - x') = k(\tau)$ .
- ▶ All *distance* kernels,  $k = k(\|x - x'\|)$  are examples of stationary kernels.
- ▶ The RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x - x'\|^2}{2\ell^2})$  is a stationary kernel. The polynomial kernel  $k_{\text{POL}}(x, x') = (x^T x + \sigma_0^2)^p$  is an example of a non-stationary kernel.
- ▶ Stationarity provides a useful *inductive bias*.

# Bochner's Theorem

## Theorem

(Bochner) A complex-valued function  $k$  on  $\mathbb{R}^P$  is the covariance function of a weakly stationary mean square continuous complex-valued random process on  $\mathbb{R}^P$  if and only if it can be represented as

$$k(\tau) = \int_{\mathbb{R}^P} e^{2\pi i s^T \tau} \psi(ds) , \quad (28)$$

where  $\psi$  is a positive finite measure.

If  $\psi$  has a density  $S(s)$ , then  $S$  is called the *spectral density* or *power spectrum* of  $k$ , and  $k$  and  $S$  are Fourier duals:

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds , \quad (29)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau . \quad (30)$$

# Review: Linear Basis Function Models

## Model Specification

$$f(x, \mathbf{w}) = \mathbf{w}^T \phi(x) \quad (31)$$

$$p(\mathbf{w}) = \mathcal{N}(0, \Sigma_w) \quad (32)$$

## Moments of Induced Distribution over Functions

$$\mathbb{E}[f(x, \mathbf{w})] = m(x) = \mathbb{E}[\mathbf{w}^T] \phi(x) = 0 \quad (33)$$

$$\text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = \mathbb{E}[f(x_i)f(x_j)] - \mathbb{E}[f(x_i)]\mathbb{E}[f(x_j)] \quad (34)$$

$$= \phi(x_i)^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi(x_j) - 0 \quad (35)$$

$$= \phi(x_i)^T \Sigma_w \phi(x_j) \quad (36)$$

- ▶  $f(x, \mathbf{w})$  is a Gaussian process,  $f(x) \sim \mathcal{N}(m, k)$  with mean function  $m(x) = 0$  and covariance kernel  $k(x_i, x_j) = \phi(x_i)^T \Sigma_w \phi(x_j)$ .
- ▶ The entire basis function model of Eqs. (31) and (32) is encapsulated as a distribution over functions with kernel  $k(x, x')$ .

# Deriving the RBF Kernel

- ▶ Start with the basis model

$$f(x) = \sum_{i=1}^J w_i \phi_i(x) , \quad (37)$$

$$w_i \sim \mathcal{N} \left( 0, \frac{\sigma^2}{J} \right) , \quad (38)$$

$$\phi_i(x) = \exp \left( -\frac{(x - c_i)^2}{2\ell^2} \right) . \quad (39)$$

- ▶ Equations (37)-(39) define a radial basis function regression model, with radial basis functions centred at the points  $c_i$ .
- ▶ Using our result for the kernel of a generalised linear model,

$$k(x, x') = \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') . \quad (40)$$

# Deriving the RBF Kernel

$$f(x) = \sum_{i=1}^J w_i \phi_i(x), \quad w_i \sim \mathcal{N}\left(0, \frac{\sigma^2}{J}\right), \quad \phi_i(x) = \exp\left(-\frac{(x - c_i)^2}{2\ell^2}\right) \quad (41)$$

$$\therefore k(x, x') = \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') \quad (42)$$

- ▶ Letting  $c_{i+1} - c_i = \Delta c = \frac{1}{J}$ , and  $J \rightarrow \infty$ , the kernel in Eq. (42) becomes a Riemann sum:

$$k(x, x') = \lim_{J \rightarrow \infty} \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') = \int_{c_0}^{c_\infty} \phi_c(x) \phi_c(x') dc \quad (43)$$

- ▶ By setting  $c_0 = -\infty$  and  $c_\infty = \infty$ , we spread the infinitely many basis functions across the whole real line, each a distance  $\Delta c \rightarrow 0$  apart:

$$k(x, x') = \int_{-\infty}^{\infty} \exp\left(-\frac{x-c}{2\ell^2}\right) \exp\left(-\frac{x'-c}{2\ell^2}\right) dc \quad (44)$$

$$= \sqrt{\pi} \ell \sigma^2 \exp\left(-\frac{(x-x')^2}{2(\ell \sqrt{\pi})^2}\right). \quad (45)$$

# Deriving the RBF Kernel

- ▶ It is remarkable we can work with infinitely many basis functions with finite amounts of computation using the *kernel trick* – replacing inner products of basis functions with kernels.
- ▶ The RBF kernel, also known as the Gaussian or squared exponential kernel, is by far the most popular kernel.
$$k_{\text{RBF}}(x, x') = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right).$$
- ▶ Recall Bochner's theorem. If we take the Fourier transform of the RBF kernel we recover a Gaussian spectral density,
$$S(s) = (2\pi\ell^2)^{D/2} \exp(-2\pi^2\ell^2 s^2)$$
 for  $x \in \mathbb{R}^D$ . Therefore the RBF kernel does not have much support for high frequency functions, since a Gaussian does not have heavy tails.
- ▶ Functions drawn from a GP with an RBF kernel are infinitely differentiable. For this reason, the RBF kernel is accused of being overly smooth and unrealistic. Nonetheless it has nice theoretical properties...

# The RBF Kernel

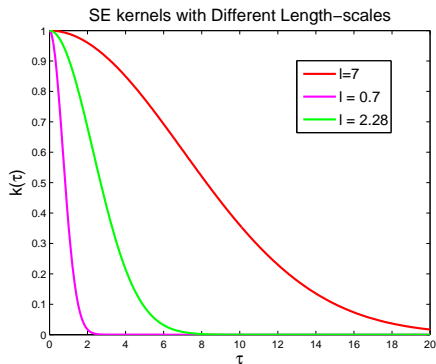


Figure: SE kernels with different length-scales, as a function of  $\tau = x - x'$ .

# Representer Theorem

A decision function  $f(x)$  can be written as

$$f(x) = \langle \mathbf{w}, \phi(x) \rangle = \left\langle \sum_{i=1}^N \alpha_i \phi(x_i), \phi(x) \right\rangle = \sum_{i=1}^N \alpha_i k(x_i, x). \quad (46)$$

- ▶ Representer theorem says this function exists with finitely many coefficients  $\alpha_i$  even when  $\phi$  is infinite dimensional (an infinite number of basis functions).
- ▶ Initially viewed as a strength of kernel methods, for datasets not exceeding e.g. ten thousand points.
- ▶ Unfortunately, the number of nonzero  $\alpha_i$  often grows linearly in the size of the training set  $N$ .
- ▶ Example: In GP regression, the predictive mean is

$$\mathbb{E}[f_* | \mathbf{y}, X, x_*] = \mathbf{k}_*^T (K + \sigma^2 I)^{-1} \mathbf{y} = \sum_{i=1}^N \alpha_i k(x_i, x_*), \quad (47)$$

where  $\alpha_i = (K + \sigma^2 I)^{-1} \mathbf{y}$ .

# Polynomial Kernel

We have already shown that the simple linear model

$$f(x, w) = \mathbf{w}^T x + b, \quad (48)$$

$$p(w) = \mathcal{N}(0, \alpha^2 I), \quad (49)$$

$$p(b) = \mathcal{N}(0, \beta^2), \quad (50)$$

corresponds to a Gaussian process with kernel

$$k_{\text{LIN}}(x, x') = \alpha^2 x^T x + \beta^2. \quad (51)$$

Samples from a GP with  $k_{\text{LIN}}(x, x')$  will thus be straight lines.

Recall that the product of two kernels is a valid kernel. The product of two linear kernels is a quadratic kernel, which gives rise to quadratic functions:

$$k_{\text{QUAD}}(x, x') = k_{\text{LIN}}(x, x') k_{\text{LIN}}(x, x'). \quad (52)$$

For example, if  $\beta = 0$ ,  $\alpha = 1$ , and  $x \in \mathbb{R}^2$ , then  $k_{\text{QUAD}}(x, x') = \phi(x)^T \phi(x')$  with  $\phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^T$ , where  $x = (x_1, x_2)$ . We can generalize to the polynomial kernel

$$k_{\text{POL}}(x, x') = (\alpha^2 x^T x + \beta^2)^p. \quad (53)$$

# The Rational Quadratic Kernel

- ▶ What if we want data varying at multiple scales?

# The Rational Quadratic Kernel

Try a scale mixture of RBF kernels. Let  $r = ||x - x'||$ .

- ▶  $k(r) = \int \exp(-\frac{r^2}{2\ell^2})p(\ell)d\ell$ .

For example, we can consider a Gamma density for  $p(\ell)$ . Letting  $\gamma = \ell^{-2}$ ,  $g(\gamma|\alpha, \beta) \propto \gamma^{\alpha-1} \exp(-\alpha\gamma/\beta)$ , with  $\beta^{-1} = \ell'^2$ , the rational quadratic (RQ) kernel is derived as

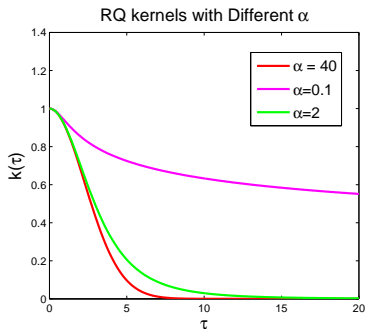
$$k_{\text{RQ}}(r) = \int_0^\infty k_{\text{RBF}}(r|\gamma)g(\gamma|\alpha, \beta)d\gamma = (1 + \frac{r^2}{2\alpha\ell'^2})^{-\alpha}. \quad (54)$$

- ▶ One could derive other interesting covariance functions using different (non-Gamma) functions for  $p(\ell)$ .

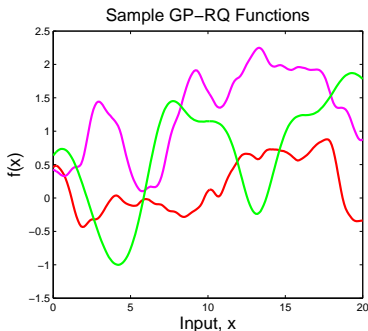
# The Rational Quadratic Kernel

$$k_{\text{RQ}}(r) = \left(1 + \frac{r^2}{2\alpha\ell^2}\right)^{-\alpha} \quad (55)$$

$$r = \|\tau\| = \|x - x'\|. \quad (56)$$



(a)



(b)

# Neural Network Kernel

- ▶ The neural network kernel (Neal, 1996) is famous for triggering research on Gaussian processes in the machine learning community.

Consider a neural network with one hidden layer:

$$f(x) = b + \sum_{i=1}^J v_i h(x; \mathbf{u}_i) . \quad (57)$$

- ▶  $b$  is a bias,  $v_i$  are the hidden to output weights,  $h$  is any bounded hidden unit transfer function,  $\mathbf{u}_i$  are the input to hidden weights, and  $J$  is the number of hidden units. Let  $b$  and  $v_i$  be independent with zero mean and variances  $\sigma_b^2$  and  $\sigma_v^2/J$ , respectively, and let the  $\mathbf{u}_i$  have independent identical distributions.

Collecting all free parameters into the weight vector  $\mathbf{w}$ ,

$$\mathbb{E}_{\mathbf{w}}[f(x)] = 0 , \quad (58)$$

$$\begin{aligned} \text{cov}[f(x), f(x')] &= \mathbb{E}_{\mathbf{w}}[f(x)f(x')] = \sigma_b^2 + \frac{1}{J} \sum_{i=1}^J \sigma_v^2 \mathbb{E}_{\mathbf{u}}[h_i(x; \mathbf{u}_i) h_i(x'; \mathbf{u}_i)] , \\ & \quad (59) \end{aligned}$$

$$= \sigma_b^2 + \sigma_v^2 \mathbb{E}_{\mathbf{u}}[h(x; \mathbf{u}) h(x'; \mathbf{u})] . \quad (60)$$

# Neural Network Kernel

$$f(x) = b + \sum_{i=1}^J v_i h(x; \mathbf{u}_i). \quad (61)$$

- ▶ Let  $h(x; \mathbf{u}) = \text{erf}(u_0 + \sum_{j=1}^P u_j x_j)$ , where  $\text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$
- ▶ Choose  $\mathbf{u} \sim \mathcal{N}(0, \Sigma)$

Then we obtain

$$k_{\text{NN}}(x, x') = \frac{2}{\pi} \sin\left(\frac{2\tilde{x}^T \Sigma \tilde{x}'}{\sqrt{(1 + 2\tilde{x}^T \Sigma \tilde{x})(1 + 2\tilde{x}'^T \Sigma \tilde{x}')}}\right), \quad (62)$$

where  $x \in \mathbb{R}^P$  and  $\tilde{x} = (1, x^T)^T$ .

# Neural Network Kernel

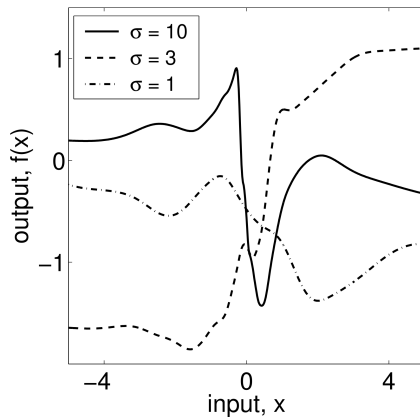


Figure: Draws from a GP with a Neural Network Kernel with Varying  $\sigma$

Rasmussen and Williams (2006)

Recall the RBF kernel

$$k_{\text{RBF}}(x, x') = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right). \quad (63)$$

- What if we want to make the length-scale of  $\ell$  input dependent, so that the resulting function is biased to vary more quickly in parts of the input space than in others?

Recall the RBF kernel

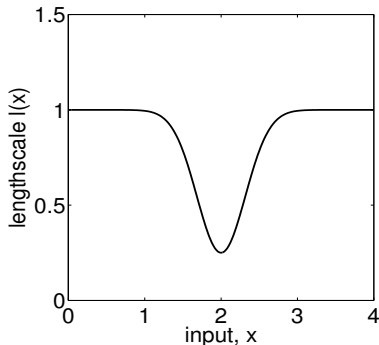
$$k_{\text{RBF}}(x, x') = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right). \quad (64)$$

- ▶ What if we want to make the length-scale of  $\ell$  input dependent, so that the resulting function is biased to vary more quickly in parts of the input space than in others?
- ▶ Just letting  $\ell \rightarrow \ell(x)$  doesn't produce a valid kernel.

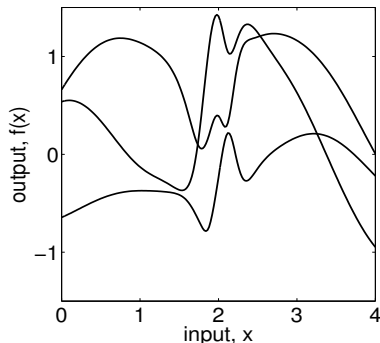
# Gibbs Kernel

$$k_{\text{Gibbs}}(x, x') = \prod_{p=1}^P \left( \frac{2l_p(x)l_p(x')}{l_p^2(x) + l_p^2(x')} \right)^{1/2} \exp \left( - \sum_{p=1}^P \frac{(x_p - x'_p)^2}{l_p^2(x) + l_p^2(x')} \right), \quad (65)$$

where  $x_p$  is the  $p^{\text{th}}$  component of  $x$ .



(a)



(b)

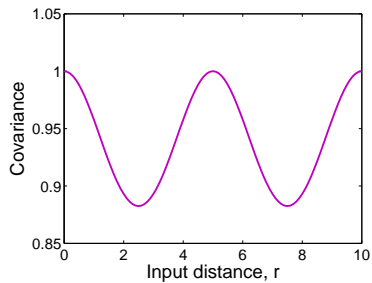
# Periodic Kernel

- ▶ Transform the inputs through a vector-valued function:  
 $\mathbf{u}(x) = (\cos(x), \sin(x))$ .
- ▶ Apply the RBF kernel in  $\mathbf{u}$  space:  $k_{\text{RBF}}(x, x') \rightarrow k_{\text{RBF}}(\mathbf{u}(x), \mathbf{u}(x'))$ .
- ▶ Recover the periodic kernel

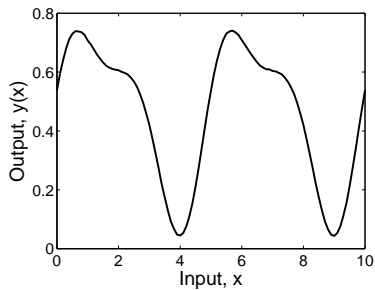
$$k_{\text{PER}}(x, x') = \exp\left(-\frac{2 \sin^2\left(\frac{x-x'}{2}\right)}{\ell^2}\right). \quad (66)$$

- ▶ Can you see anything unusual about this kernel?

# Periodic Kernel



(a)



(b)

# Non-Stationary Kernels

- ▶ A stationary kernel is invariant to translations of the input space:  
 $k = k(\tau), \tau = x - x'$ .
- ▶ Intuitively, this means the properties of the function are similar across different regions of the input domain.
- ▶ How might we make other non-stationary kernels, besides the Gibbs kernel?

# Non-Stationary Kernels

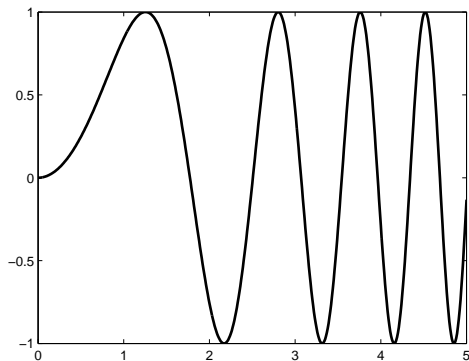
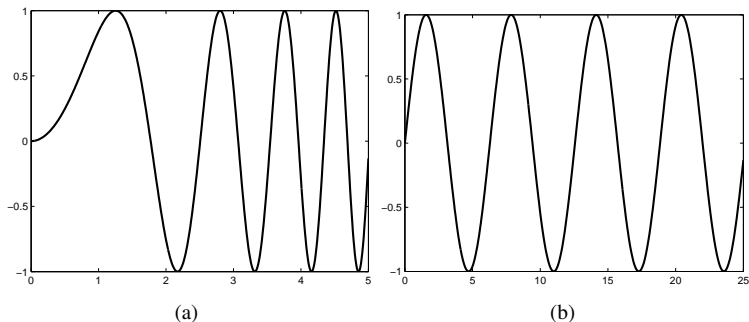


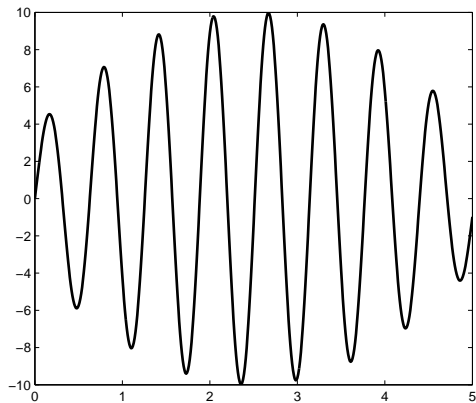
Figure: Non-stationary function

# Non-Stationary Kernels



**Figure:** Warp the inputs (in this case,  $x \rightarrow x^2$  to go from non-stationary function to a stationary function). E.g., apply  $k(g(x), g(x'))$  to the data, where  $g$  is a warping function.

# Non-Stationary Kernels



# Non-Stationary Kernels

- ▶ Warp the input space:  $k(x, x') \rightarrow k(g(x), g(x'))$  where  $g$  is an arbitrary warping function.
- ▶ Modulate the amplitude of the kernel. If  $f(x) \sim \mathcal{GP}(0, k(x, x'))$  then  $a(x)f(x)$  has kernel  $a(x)k(x, x')a(x')$ , conditioned on  $a(x)$ .
- ▶ What would happen if we tried  $w_1(x)f_1(x) + w_2(x)f_2(x)$  where  $f_1$  and  $f_2$  are GPs with different kernels?
- ▶ How about  $\sigma(w_1(x))f_1(x) + (1 - \sigma(w_1(x)))f_2(x)$ ?

- ▶ The RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x-x'\|^2}{2\ell^2})$  is criticized for being too smooth.
- ▶ How might we create a drop-in replacement, while retaining useful inductive biases?

# Matérn Kernel

- ▶ The RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x-x'\|^2}{2\ell^2})$  is criticized for being too smooth.
- ▶ How might we create a drop-in replacement, while retaining useful inductive biases?
- ▶ Could replace the Euclidean distance measure with an absolute distance measure... then we recover the Ornstein-Uhlenbeck kernel:  
 $k_{\text{OU}}(x, x') = \exp(\|x - x'\|/\ell)$ . The velocity of a particle undergoing brownian motion is described by a GP with the OU kernel.

# Matérn Kernel

- ▶ The RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x-x'\|^2}{2\ell^2})$  is criticized for being too smooth.
- ▶ How might we create a drop-in replacement, while retaining useful inductive biases?
- ▶ Could replace the Euclidean distance measure with an absolute distance measure... then we recover the Ornstein-Uhlenbeck kernel:  
 $k_{\text{OU}}(x, x') = \exp(\|x - x'\|/\ell)$ . The velocity of a particle undergoing brownian motion is described by a GP with the OU kernel.
- ▶ Recall that stationary kernels  $k(\tau)$ ,  $\tau = x - x'$  and spectral densities are Fourier duals of one another:

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds, \quad (67)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau. \quad (68)$$

If we take the Fourier transform of the RBF kernel, we recover a Gaussian spectral density... But we can go from spectral densities to kernels too...

# Matérn Kernel

- ▶ The RBF kernel  $k_{\text{RBF}}(x, x') = a^2 \exp(-\frac{\|x-x'\|^2}{2\ell^2})$  is criticized for being too smooth.
- ▶ How might we create a drop-in replacement, while retaining useful inductive biases?
- ▶ Could replace the Euclidean distance measure with an absolute distance measure... then we recover the Ornstein-Uhlenbeck kernel:  
 $k_{\text{OU}}(x, x') = \exp(\|x - x'\|/\ell)$ . The velocity of a particle undergoing brownian motion is described by a GP with the OU kernel.
- ▶ Recall that stationary kernels  $k(\tau)$ ,  $\tau = x - x'$  and spectral densities are Fourier duals of one another:

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds, \quad (69)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau. \quad (70)$$

If we take the Fourier transform of the RBF kernel, we recover a Gaussian spectral density... But we can go from spectral densities to kernels too...

- ▶ If we use a Student- $t$  spectral density for  $S(s)$ , and take the inverse Fourier transform, we recover the *Matérn* kernel.

# Matérn Kernel



$$k_{\text{Matérn}}(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}|x - x'|}{\ell} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}|x - x'|}{\ell} \right), \quad (71)$$

where  $K_\nu$  is a modified Bessel function.

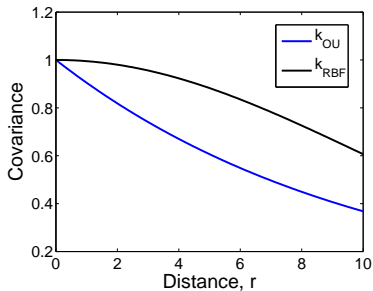
- ▶ In one dimension, and when  $\nu + 1/2 = p$ , for some natural number  $p$ , the corresponding GP is a continuous time AR( $p$ ) process.
- ▶ By setting  $\nu = 1$ , we obtain the *Ornstein-Uhlenbeck* (OU) kernel,

$$k_{\text{OU}}(x, x') = \exp\left(-\frac{\|x - x'\|}{\ell}\right). \quad (72)$$

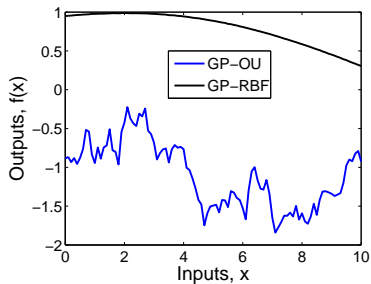
- ▶ The Matérn kernel does not have *concentration of measure* problems for high dimensional inputs to the extent of the RBF (Gaussian) kernel (Fastfood: Le, Sarlos, Smola, ICML 2013).
- ▶ The kernel gives rise to a Markovian process (and classical filtering and smoothing algorithms can be applied).

# Matérn Kernel

OU and RBF kernels both with lengthscale  $\ell = 10$ .

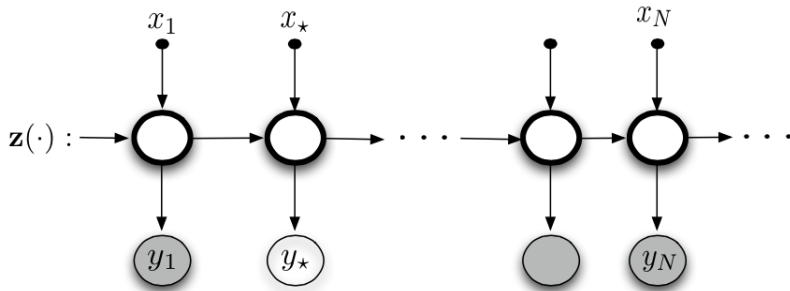


(a)



(b)

# Matérn Kernel



From Yunus Saatchi's PhD thesis, *Scalable Inference for Structured Gaussian Process Models*, 2011.

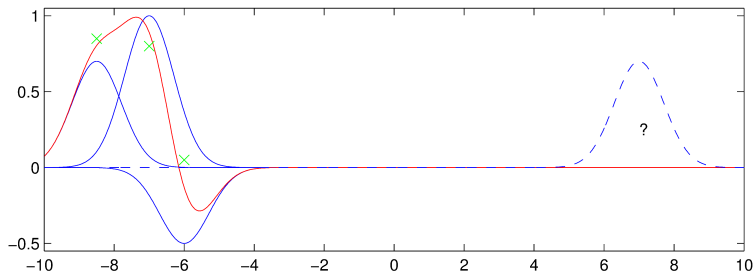
- ▶ Are Gaussian processes Bayesian nonparametric models?

# Nonparametric Kernels

- ▶ For a Gaussian process  $f(x)$  to be non-parametric,  $f(x_i)|\mathbf{f}_{-i}$ , where  $\mathbf{f}_{-i}$  is any collection of function values excluding  $f(x_i)$ , must be free to take any value in  $\mathbb{R}$ .
- ▶ For this freedom to be possible it is a necessary (but not sufficient) condition for the kernel of the Gaussian process to be derived from an infinite basis function expansion.
- ▶ Nonparametric kernels allow for a great amount of flexibility: the amount of information the model can represent grows with the amount of available data.

# Nonparametric RBF vs Finite Dimensional Analogue

- The parametric analogue to a GP with a non-parametric RBF kernel becomes *more* confident in its predictions, the further away we get from the data!



Rasmussen, MLSS Cambridge, 2009.

# Simple Random Walk

- ▶ Discrete time auto-regressive model



$$f(t) = a f(t-1) + \epsilon(t), \quad (73)$$

$$\epsilon(t) \sim \mathcal{N}(0, 1), \quad (74)$$

$$a \in \mathbb{R}, \quad (75)$$

$$t = 1, 2, 3, 4, \dots \quad (76)$$

$$(77)$$

- ▶ Is this model a Gaussian process?

# Gaussian Process Covariance Kernels

Let  $\tau = x - x'$ :

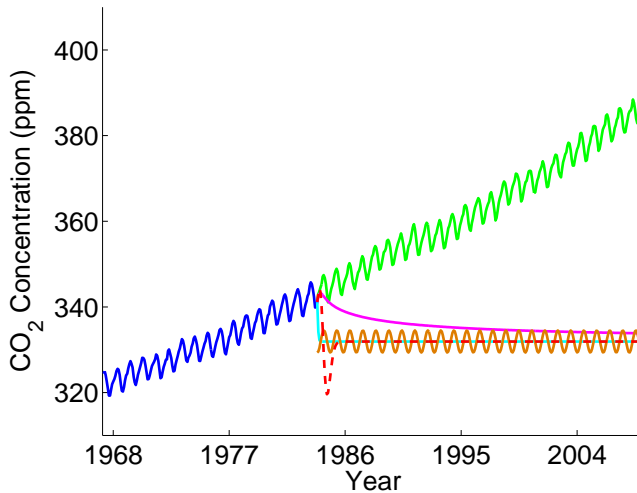
$$k_{\text{SE}}(\tau) = \exp(-0.5\tau^2/\ell^2) \quad (78)$$

$$k_{\text{MA}}(\tau) = a(1 + \frac{\sqrt{3}\tau}{\ell}) \exp(-\frac{\sqrt{3}\tau}{\ell}) \quad (79)$$

$$k_{\text{RQ}}(\tau) = (1 + \frac{\tau^2}{2\alpha\ell^2})^{-\alpha} \quad (80)$$

$$k_{\text{PE}}(\tau) = \exp(-2 \sin^2(\pi \tau \omega)/\ell^2) \quad (81)$$

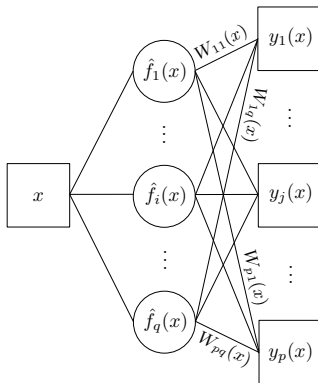
# CO<sub>2</sub> Extrapolation with Standard Kernels



“How can Gaussian processes possibly replace neural networks? Did we throw the baby out with the bathwater?”

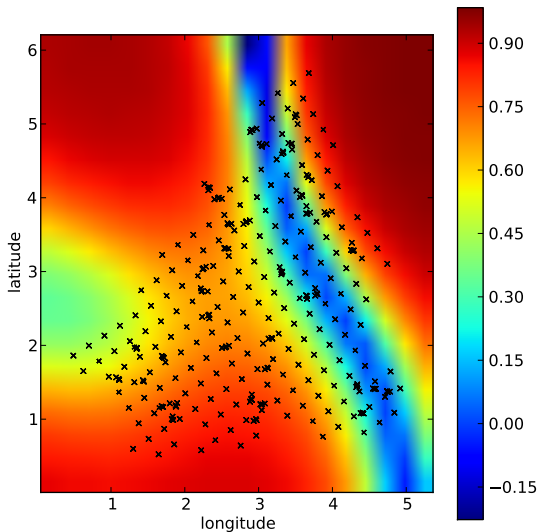
David MacKay, 1998.

# More Expressive Covariance Functions



*Gaussian Process Regression Networks.* Wilson et. al, ICML 2012.

# Gaussian Process Regression Network



# Expressive Covariance Functions

- ▶ GPs in Bayesian neural network like architectures. (Salakhutdinov and Hinton, 2008; Wilson et. al, 2012; Damianou and Lawrence, 2012).  
Task specific, difficult inference, no closed form kernels.
- ▶ Compositions of kernels. (Archambeau and Bach, 2011; Durrande et. al, 2011; Rasmussen and Williams, 2006).  
In the general case, difficult to interpret, difficult inference, struggle with over-fitting.

Can learn almost nothing about the covariance function of a stochastic process from a single realization, if we assume that the covariance function could be *any* positive definite function. Most commonly one assumes a restriction to *stationary* kernels, meaning that covariances are invariant to translations in the input space.

# Bochner's Theorem

## Theorem

(Bochner) A complex-valued function  $k$  on  $\mathbb{R}^P$  is the covariance function of a weakly stationary mean square continuous complex-valued random process on  $\mathbb{R}^P$  if and only if it can be represented as

$$k(\tau) = \int_{\mathbb{R}^P} e^{2\pi i s^T \tau} \psi(ds) , \quad (82)$$

where  $\psi$  is a positive finite measure.

If  $\psi$  has a density  $S(s)$ , then  $S$  is called the *spectral density* or *power spectrum* of  $k$ , and  $k$  and  $S$  are Fourier duals:

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds , \quad (83)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau . \quad (84)$$

$k$  and  $S$  are Fourier duals:

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds, \quad (85)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau. \quad (86)$$

- ▶ If we can approximate  $S(s)$  to arbitrary accuracy, then we can approximate any stationary kernel to arbitrary accuracy.
- ▶ We can model  $S(s)$  to arbitrary accuracy, since scale-location mixtures of Gaussians can approximate any distribution to arbitrary accuracy.
- ▶ A scale-location mixture of Gaussians can flexibly model many distributions, and thus many covariance kernels, even with a small number of components.

# Kernels for Pattern Discovery

Let  $\tau = x - x' \in \mathbb{R}^P$ . From Bochner's Theorem,

$$k(\tau) = \int_{\mathbb{R}^P} S(s) e^{2\pi i s^T \tau} ds \quad (87)$$

For simplicity, assume  $\tau \in \mathbb{R}^1$  and let

$$S(s) = [\mathcal{N}(s; \mu, \sigma^2) + \mathcal{N}(-s; \mu, \sigma^2)]/2. \quad (88)$$

Then

$$k(\tau) = \exp\{-2\pi^2 \tau^2 \sigma^2\} \cos(2\pi \tau \mu). \quad (89)$$

More generally, if  $S(s)$  is a symmetrized mixture of diagonal covariance Gaussians on  $\mathbb{R}^P$ , with covariance matrix  $\mathbf{M}_q = \text{diag}(v_q^{(1)}, \dots, v_q^{(P)})$ , then

$$k(\tau) = \sum_{q=1}^Q w_q \cos(2\pi \tau_p \mu_q^{(p)}) \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\}. \quad (90)$$

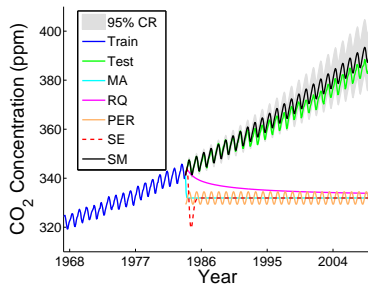
# GP Model for Pattern Extrapolation

- ▶ Observations  $y(x) \sim \mathcal{N}(y(x); f(x), \sigma^2)$  (can easily be relaxed).
- ▶  $f(x) \sim \mathcal{GP}(0, k_{\text{SM}}(x, x' | \theta))$  ( $f(x)$  is a GP with SM kernel).
- ▶  $k_{\text{SM}}(x, x' | \theta)$  can approximate many different kernels with different settings of its hyperparameters  $\theta$ .
- ▶ *Learning* involves training these hyperparameters through maximum marginal likelihood optimization (using BFGS)

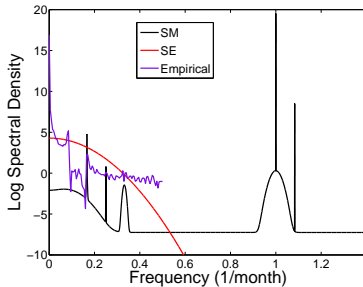
$$\log p(\mathbf{y} | \theta, X) = \underbrace{-\frac{1}{2} \mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2} \log |K_\theta + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2} \log(2\pi). \quad (91)$$

- ▶ Once hyperparameters are trained as  $\hat{\theta}$ , making predictions using  $p(f_* | \mathbf{y}, X_*, \hat{\theta})$ , which can be expressed in closed form.

# Results, CO<sub>2</sub>

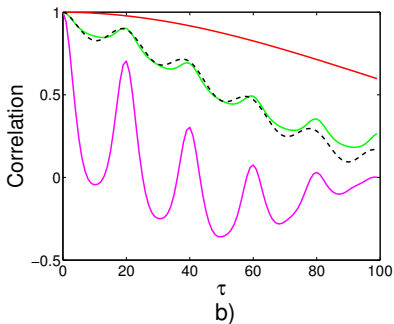
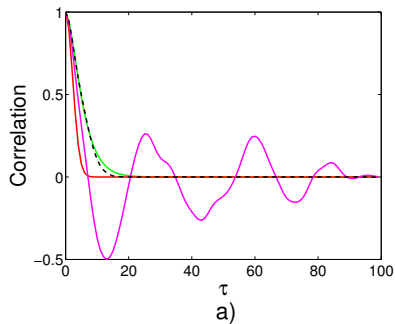


(c)

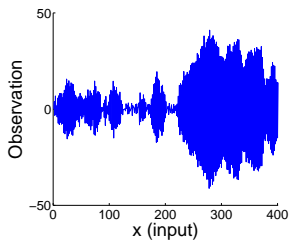


(d)

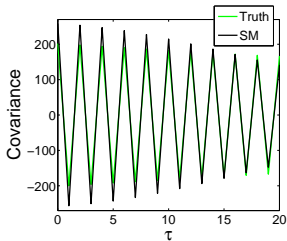
# Results, Reconstructing Standard Covariances



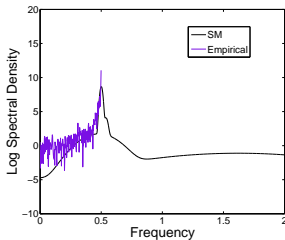
# Results, Negative Covariances



(e)

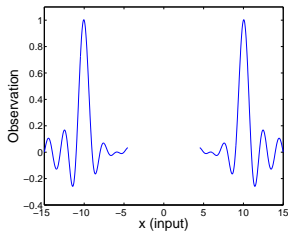


(f)

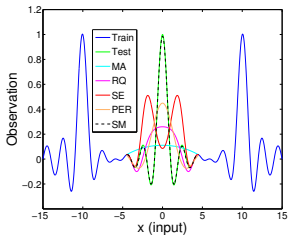


(g)

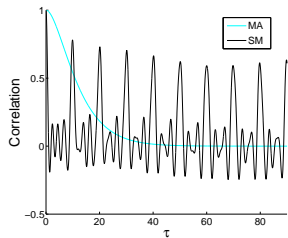
# Results, Sinc Pattern



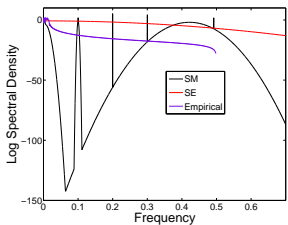
(h)



(i)

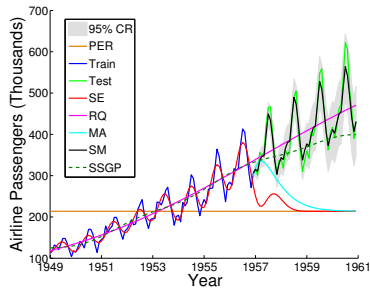


(j)

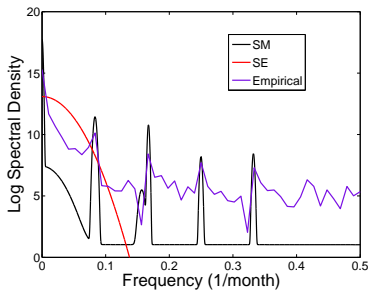


(k)

# Results, Airline Passengers



(l)



(m)

# Scaling up kernel machines

- ▶ Expressive kernels will be most valuable on large datasets.
- ▶ Computational bottlenecks for GPs:
  - ▶ Inference:  $(K_\theta + \sigma^2 I)^{-1} \mathbf{y}$  for  $n \times n$  matrix  $K$ .
  - ▶ Learning:  $\log |K_\theta + \sigma^2 I|$ , for marginal likelihood evaluations needed to learn  $\theta$ .
- ▶ Both inference and learning naively require  $\mathcal{O}(n^3)$  operations and  $\mathcal{O}(n^2)$  storage (typically from computing a Cholesky decomposition of  $K$ ). Afterwards, the predictive mean and variance cost  $\mathcal{O}(n)$  and  $\mathcal{O}(n^2)$  per test point.

# Inference and Learning

1. **Learning:** Optimize marginal likelihood,

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

with respect to kernel hyperparameters  $\boldsymbol{\theta}$ .

2. **Inference:** Conditioned on kernel hyperparameters  $\boldsymbol{\theta}$ , form the predictive distribution for test inputs  $X_*$ :

$$\mathbf{f}_*|X_*, X, \mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)) ,$$

$$\bar{\mathbf{f}}_* = K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}\mathbf{y} ,$$

$$\text{cov}(\mathbf{f}_*) = K_{\boldsymbol{\theta}}(X_*, X_*) - K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}K_{\boldsymbol{\theta}}(X, X_*) .$$

## Three Families of Approaches

- ▶ Approximate non-parametric kernels in a finite basis ‘dual space’. Requires  $\mathcal{O}(m^2n)$  computations and  $\mathcal{O}(m)$  storage for  $m$  basis functions. Examples: SSGP, Random Kitchen Sinks, Fastfood, À la Carte.
- ▶ Inducing point based sparse approximations. Examples: SoR, FITC, KISS-GP.
- ▶ Exploit existing structure in  $K$  to quickly (and exactly) solve linear systems and log determinants. Examples: Toeplitz and Kronecker methods.

# Parametric Expansions via Random Basis Functions

- ▶ Return to Bochner's Theorem

$$k(\tau) = \int S(s) e^{2\pi i s^T \tau} ds, \quad (92)$$

$$S(s) = \int k(\tau) e^{-2\pi i s^T \tau} d\tau. \quad (93)$$

- ▶ We can treat  $S(s)$  as a probability distribution and sample from it, to approximate the integral for  $k(\tau)$ !
- ▶ It is a valid Monte Carlo procedure to sample the pairs  $\{s_j, -s_j\}$  from  $S(s)$ :

$$k(\tau) \approx \frac{1}{2J} \sum_{j=1}^J [\exp(2\pi i s_j^T \tau) + \exp(-2\pi i s_j^T \tau)], \quad s_j \sim S(s) \quad (94)$$

$$= \frac{1}{J} \sum_{j=1}^J \cos(2\pi s_j^T \tau) \quad (95)$$

- ▶ This is exactly the covariance function we get if we use a linear basis function model with trigonometric basis functions! Use the basis function representation with finite  $J$  for computational efficiency.

# Scaling a Gaussian process: inducing inputs

- ▶ Gaussian process  $f$  and  $f_*$  evaluated at  $n$  training points and  $J$  testing points.
- ▶  $m \ll n$  inducing points  $u$ ,  $p(u) = \mathcal{N}(0, K_{u,u})$
- ▶  $p(f_*, f) = \int p(f_*, f, u) du = \int p(f_*, f | u) p(u) du$
- ▶ Assume that  $f$  and  $f_*$  are conditionally independent given  $u$ :

$$p(f_*, f) \approx q(f_*, f) = \int q(f_* | u) q(f | u) p(u) du \quad (96)$$

- ▶ Exact conditional distributions

$$p(f | u) = \mathcal{N}(K_{f,u} K_{u,u}^{-1} u, K_{f,f} - Q_{f,f}) \quad (97)$$

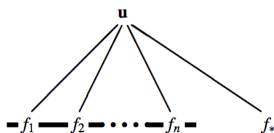
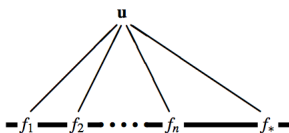
$$p(f_* | u) = \mathcal{N}(K_{f_*,u} K_{u,u}^{-1} u, K_{f_*,f_*} - Q_{f_*,f_*}) \quad (98)$$

$$Q_{a,b} = K_{a,u} K_{u,u}^{-1} K_{u,b} \quad (99)$$

- ▶ Cost for predictions reduced from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(m^2 n)$  where  $m \ll n$ .
- ▶ Different inducing approaches correspond to different additional assumptions about  $q(f | u)$  and  $q(f_* | u)$ .

# Inducing Point Methods

The inducing points act as a communication channel between the GP evaluated at the training and test points,  $\mathbf{f}$  and  $\mathbf{f}_*$ :



# Subset of Regression (SoR)

The subset of regressors method uses deterministic conditional distributions with exact means:

$$q(\mathbf{f}|\mathbf{u}) = \mathcal{N}(K_{X,U}K_{U,U}^{-1}\mathbf{u}, \mathbf{0}) \quad (100)$$

$$q(\mathbf{f}_*|\mathbf{u}) = \mathcal{N}(K_{X_*,U}K_{U,U}^{-1}\mathbf{u}, \mathbf{0}) \quad (101)$$

$$(102)$$

Integrate away  $\mathbf{u}$  via

$$p(\mathbf{f}_*, \mathbf{f}) \approx q(\mathbf{f}_*, \mathbf{f}) = \int q(\mathbf{f}_*|\mathbf{u})q(\mathbf{f}|\mathbf{u})p(\mathbf{u})d\mathbf{u} \quad (103)$$

to obtain the joint distribution

$$q_{\text{SoR}} = \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} Q_{X,X} & Q_{X,X_*} \\ Q_{X_*,X} & Q_{X_*,X_*} \end{bmatrix}\right). \quad (104)$$

$$Q_{a,b} = K_{a,u}K_{u,u}^{-1}K_{u,b} \quad (105)$$

# Subsets of Regressors

The predictive conditional can then be derived as before:

$$q_{\text{SoR}}(\mathbf{f}_* | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}, A) \quad (106)$$

$$\boldsymbol{\mu} = \mathbf{Q}_{X_*, X}(\mathbf{Q}_{X, X} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \quad (107)$$

$$A = \mathbf{Q}_{X_*, X_*} - \mathbf{Q}_{X_*, X}(\mathbf{Q}_{X, X} + \sigma^2)^{-1} \mathbf{Q}_{X, X_*} \quad (108)$$

This method can be viewed as replacing the exact covariance function  $k$  with an approximate covariance function

$$k_{\text{SoR}}(x_i, x_j) = k(x_i, U) K_{U, U}^{-1} k(U, x_j) \quad (109)$$

which admits fast computations.

# Subsets of Regressors

The SoR covariance matrix is

$$\overbrace{K_{\text{SoR}}(X, X)}^{n \times n} = \overbrace{K_{X,U}}^{n \times m} \overbrace{K_{U,U}^{-1}}^{m \times m} \overbrace{K_{U,X}}^{m \times n} \quad (110)$$

- ▶ For  $m < n$ , this is a low rank covariance matrix, corresponding to a degenerate (finite basis) Gaussian process.
- ▶ As a result, for  $n$  large, SoR tends to underestimate uncertainty.

FITC, the most popular inducing point method, uses the exact test conditional, and a factorized training conditional:

$$q_{\text{FITC}}(\mathbf{f}|\mathbf{u}) = \prod_{i=1}^n p(f_i|\mathbf{u}) \quad (111)$$

$$q_{\text{FITC}}(\mathbf{f}_*|\mathbf{u}) = p(\mathbf{f}_*|\mathbf{u}) . \quad (112)$$

Integrating away  $\mathbf{u}$ , we can derive the FITC approximate kernel as:

$$\tilde{k}_{\text{SoR}}(x, z) = K_{x,U} K_{U,U}^{-1} K_{U,z} , \quad (113)$$

$$\tilde{k}_{\text{FITC}}(x, z) = \tilde{k}_{\text{SoR}}(x, z) + \delta_{xz} \left( k(x, z) - \tilde{k}_{\text{SoR}}(x, z) \right) . \quad (114)$$

FITC replaces the diagonal of the SoR approximation with the true diagonal of  $k$ . FITC corresponds to a non-parametric GP.

# Kronecker methods

Suppose

- ▶ If  $x \in \mathbb{R}^P$ ,  $k$  decomposes as a product of kernels across each input dimension:  $k(x_i, x_j) = \prod_{p=1}^P k^p(x_i^p, x_j^p)$  (e.g., the RBF kernel has this property).
- ▶ Suppose the inputs  $x \in \mathcal{X}$  are on a multidimensional grid  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_P \subset \mathbb{R}^P$ .

Then

- ▶  $K$  decomposes into a Kronecker product of matrices over each input dimension  $K = K^1 \otimes \dots \otimes K^P$ .
- ▶ The eigendecomposition of  $K$  into  $QVQ$  also decomposes:  $Q = Q^1 \otimes \dots \otimes Q^P$ ,  $V = Q^1 \otimes \dots \otimes Q^P$ . Assuming equal cardinality for each input dimension, we can thus eigendecompose an  $N \times N$  matrix  $K$  in  $\mathcal{O}(PN^{3/P})$  operations instead of  $\mathcal{O}(N^3)$  operations.

# Kronecker methods

Suppose

- ▶ If  $x \in \mathbb{R}^P$ ,  $k$  decomposes as a product of kernels across each input dimension:  $k(x_i, x_j) = \prod_{p=1}^P k^p(x_i^p, x_j^p)$  (e.g., the RBF kernel has this property).
- ▶ Suppose the inputs  $x \in \mathcal{X}$  are on a multidimensional grid  $\mathcal{X} = \mathcal{X}_1 \times \cdots \times \mathcal{X}_P \subset \mathbb{R}^P$ .

Then

- ▶  $K$  decomposes into a Kronecker product of matrices over each input dimension  $K = K^1 \otimes \cdots \otimes K^P$ .
- ▶ The eigendecomposition of  $K$  into  $QVQ$  also decomposes:  $Q = Q^1 \otimes \cdots \otimes Q^P$ ,  $V = Q^1 \otimes \cdots \otimes Q^P$ . Assuming equal cardinality for each input dimension, we can thus eigendecompose an  $N \times N$  matrix  $K$  in  $\mathcal{O}(PN^{3/P})$  operations instead of  $\mathcal{O}(N^3)$  operations.

Then inference and learning are highly efficient:

▶

$$(K + \sigma^2 I)^{-1} \mathbf{y} = (QVQ^T + \sigma^2 I)^{-1} \mathbf{y} = Q(V + \sigma^2 I)^{-1} Q^T \mathbf{y}, \quad (115)$$

$$\log |K + \sigma^2 I| = \log |QVQ^T + \sigma^2 I| = \sum_{i=1}^N \log(\lambda_i + \sigma^2), \quad (116)$$

# Kronecker Methods

- ▶ We assumed that the inputs  $x \in \mathcal{X}$  are on a multidimensional grid  $\mathcal{X} = \mathcal{X}_1 \times \cdots \times \mathcal{X}_p \subset \mathbb{R}^P$ .
- ▶ How might we relax this assumption, to use Kronecker methods if there are gaps (missing data) in our multidimensional grid?

# Kronecker Methods

- ▶ Assume imaginary points that complete the grid
- ▶ Place infinite noise on these points so they have no effect on inference
- ▶ The relevant matrices are no longer Kronecker, but we can get around this using pre-conditioned conjugate gradients, an iterative linear solver.

# Kronecker Methods with Missing Data

- ▶ Assuming we have a dataset of  $M$  observations which are not necessarily on a grid, we propose to form a complete grid using  $W$  imaginary observations,  $\mathbf{y}_W \sim \mathcal{N}(\mathbf{f}_W, \epsilon^{-1}I_W)$ ,  $\epsilon \rightarrow 0$ .
- ▶ The total observation vector  $\mathbf{y} = [\mathbf{y}_M, \mathbf{y}_W]^T$  has  $N = M + W$  entries:  $\mathbf{y} = \mathcal{N}(\mathbf{f}, D_N)$ , where the noise covariance matrix  $D_N = \text{diag}(D_M, \epsilon^{-1}I_W)$ ,  $D_M = \sigma^2 I_M$ .
- ▶ The imaginary observations  $\mathbf{y}_W$  have *no corrupting effect* on inference: the moments of the resulting predictive distribution are exactly the same as for the standard predictive distribution, namely  $\lim_{\epsilon \rightarrow 0} (K_N + D_N)^{-1} \mathbf{y} = (K_M + D_M)^{-1} \mathbf{y}_M$ .

# Kronecker Methods with Missing Inputs

- ▶ We use preconditioned conjugate gradients to compute  $(K_N + D_N)^{-1} \mathbf{y}$ . We use the preconditioning matrix  $C = D_N^{-1/2}$  to solve  $C^T (K_N + D_N) C \mathbf{z} = C^T \mathbf{y}$ . The preconditioning matrix  $C$  speeds up convergence by ignoring the imaginary observations  $\mathbf{y}_W$ .
- ▶ For the log complexity in the marginal likelihood (used in hyperparameter learning),

$$\log |K_M + D_M| = \sum_{i=1}^M \log(\lambda_i^M + \sigma^2) \approx \sum_{i=1}^M \log(\tilde{\lambda}_i^M + \sigma^2), \quad (117)$$

where  $\tilde{\lambda}_i^M = \frac{M}{N} \lambda_i^N$  for  $i = 1, \dots, M$ .

# Spectral Mixture Product Kernel

- ▶ The spectral mixture kernel, in its standard form, does not quite have Kronecker structure.
- ▶ Introduce a *spectral mixture product kernel*, which takes a product of across input dimensions of one dimensional spectral mixture kernels.

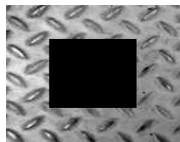
$$k_{\text{SMP}}(\tau|\boldsymbol{\theta}) = \prod_{p=1}^P k_{\text{SM}}(\tau_p|\boldsymbol{\theta}_p). \quad (118)$$

- ▶ Observations  $y(x) \sim \mathcal{N}(y(x); f(x), \sigma^2)$  (can easily be relaxed).
- ▶  $f(x) \sim \mathcal{GP}(0, k_{\text{SMP}}(x, x' | \theta))$  ( $f(x)$  is a GP with SMP kernel).
- ▶  $k_{\text{SMP}}(x, x' | \theta)$  can approximate many different kernels with different settings of its hyperparameters  $\theta$ .
- ▶ *Learning* involves training these hyperparameters through maximum marginal likelihood optimization (using BFGS)

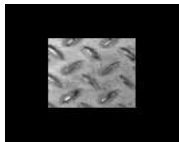
$$\log p(\mathbf{y} | \theta, X) = \overbrace{-\frac{1}{2} \mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}^{\text{model fit}} - \overbrace{\frac{1}{2} \log |K_\theta + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2} \log(2\pi). \quad (119)$$

- ▶ Once hyperparameters are trained as  $\hat{\theta}$ , making predictions using  $p(f_* | \mathbf{y}, X_*, \hat{\theta})$ , which can be expressed in closed form.
- ▶ Exploit Kronecker structure for fast exact inference and learning (and extend Kronecker methods to allow for non-grid data). *Exact* inference and learning requires  $\mathcal{O}(PN^{\frac{P+1}{P}})$  operations and  $\mathcal{O}(PN^{\frac{2}{P}})$  storage, compared to  $\mathcal{O}(N^3)$  operations and  $\mathcal{O}(N^2)$  storage, for  $N$  datapoints, and  $P$  input dimensions.

# Results



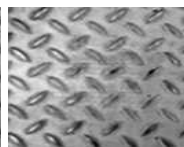
(a) Train



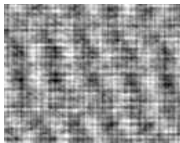
(b) Test



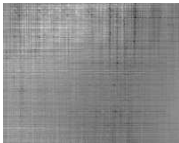
(c) Full



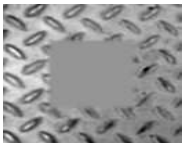
(d) GPatt



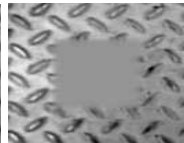
(e) SSGP



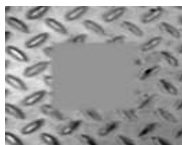
(f) FITC



(g) GP-SE



(h) GP-MA



(i) GP-RQ

# Results: Extrapolation and Interpolation with Shadows



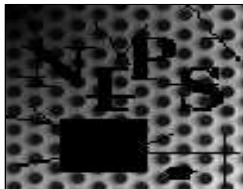
(a) Train



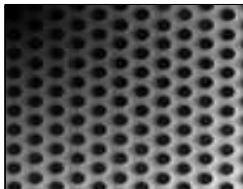
(b) GPatt



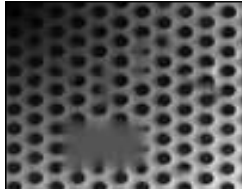
(c) GP-MA



(d) Train

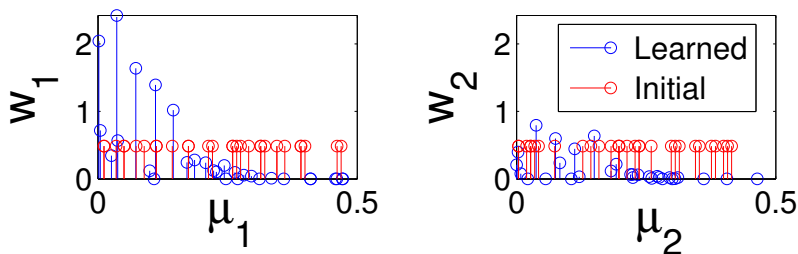


(e) GPatt



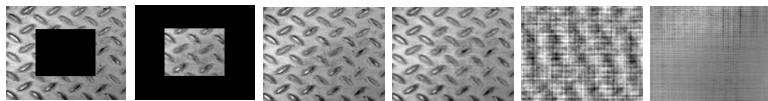
(f) GP-MA

# Automatic Model Selection via Marginal Likelihood



- ▶ Simple initialisation
- ▶ The marginal likelihood shrinks weights of extraneous components to zero through the  $\log |K|$  complexity penalty.

# Results



(a) Train

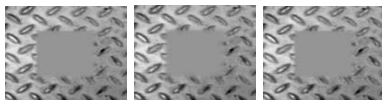
(b) Test

(c) Full

(d) GPatt

(e) SSGP

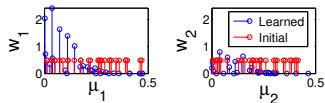
(f) FITC



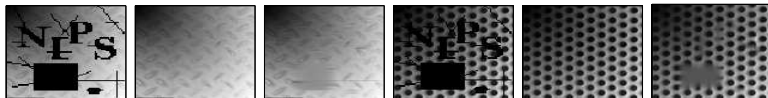
(g) GP-SE

(h) GP-MA

(i) GP-RQ



(j) GPatt Initialisation



(k) Train

(l) GPatt

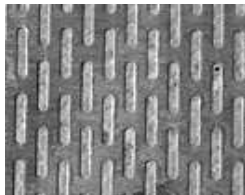
(m) GP-MA

(n) Train

(o) GPatt

(p) GP-MA

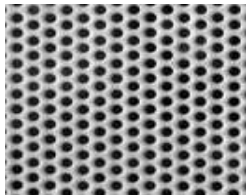
# More Patterns



(a) Rubber mat



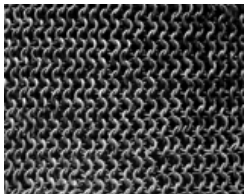
(b) Tread plate



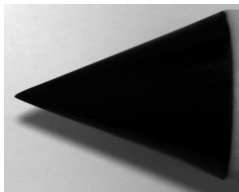
(c) Pores



(d) Wood

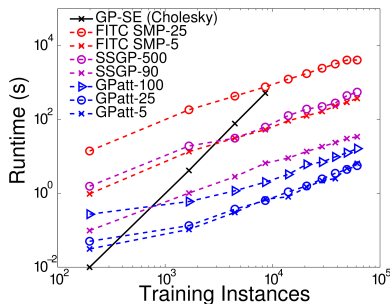


(e) Chain mail

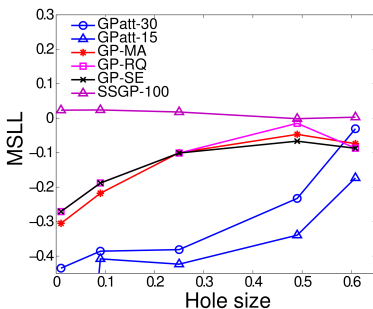


(f) Cone

# Speed and Accuracy Stress Tests

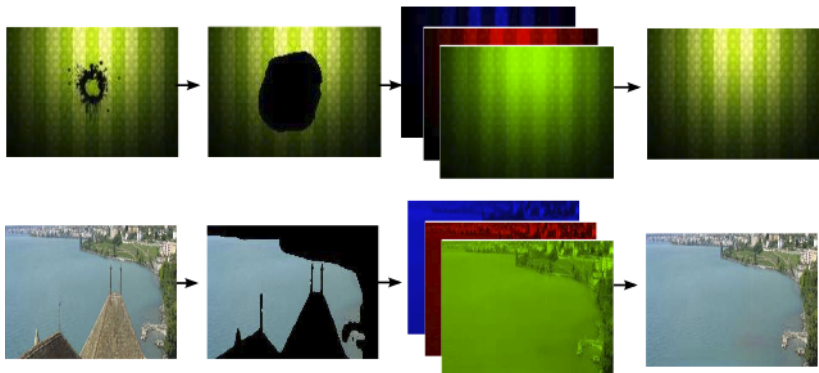


(a) Runtime Stress Test

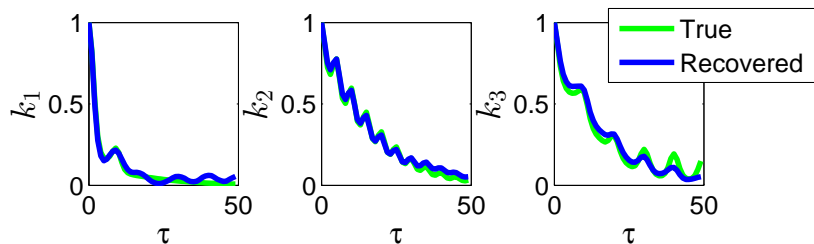


(b) Accuracy Stress Test

# Image Inpainting

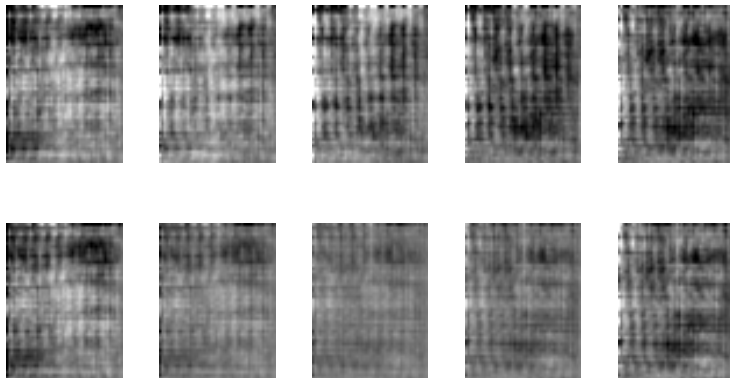


# Recovering Sophisticated Out of Class Kernels



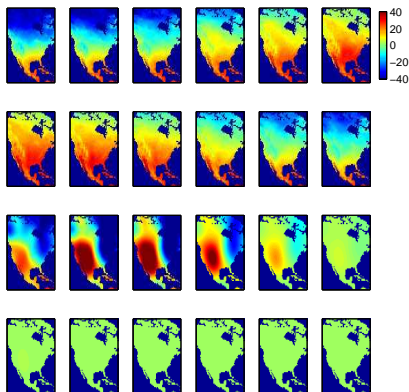
# Video Extrapolation

- ▶ GPatt makes almost no assumptions about the correlation structures across input dimensions: it can automatically discover both temporal and spatial correlations!
- ▶ Top row: True frames taken from the middle of a movie. Bottom row: Predicted sequence of frames (all are forecast together).
- ▶ 112,500 datapoints. GPatt training time is under 5 minutes.



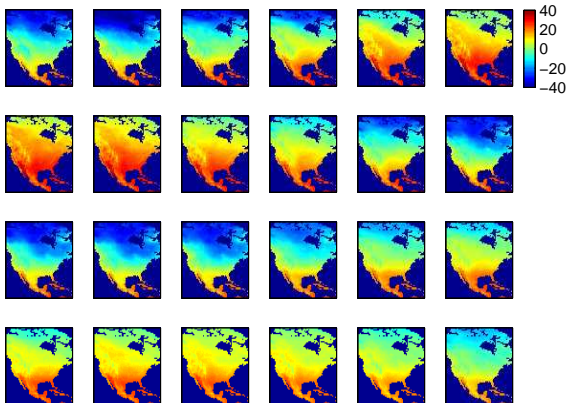
# Land Surface Temperature Forecasting

- ▶ Train using 9 years of temperature data. First two rows are the last 12 months of training data, last two rows is a 12 month ahead forecast. 300,000 data points, with 40% missing data (from ocean).
- ▶ Predictions using GP-SE (GP with an SE or RBF kernel), and Kronecker Inference.

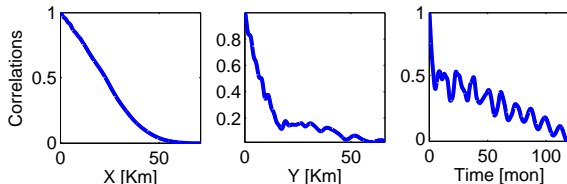


# Land Surface Temperature Forecasting

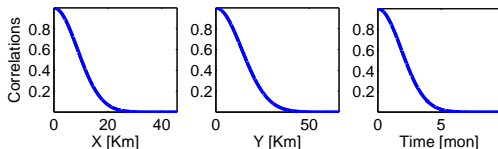
- ▶ Train using 9 years of temperature data. First two rows are the last 12 months of training data, last two rows is a 12 month ahead forecast. 300,000 data points, with 40% missing data (from ocean).
- ▶ Predictions using GPatt. Training time < 30 minutes.



# Learned Kernels for Land Surface Temperatures



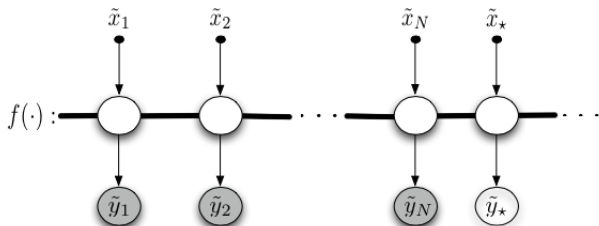
(a) Learned GPatt Kernel for Temperatures



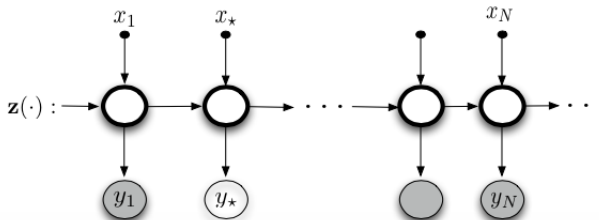
(b) Learned GP-SE Kernel for Temperatures

- The learned GPatt kernel tells us interesting properties of the data. In this case, the learned kernels are heavy tailed and quasi-periodic.

# Building Gauss-Markov Processes



Sort dataset &  
augment state space



# Generalising inducing point methods

Blackboard discussion