

# Gaussian Processes

Andrew Gordon Wilson

<https://cims.nyu.edu/~andrewgw>  
Courant Institute of Mathematical Sciences  
Center for Data Science  
New York University

Pavel's Machine Learning Class  
April 23, 2026

# The Function-Space View

*The parameters in isolation are completely divorced from the statistical properties of a model. Yet we focus most of our effort on learning parameters  $\mathbf{w}$ . What we really care about are how those parameters  $\mathbf{w}$  combine with a functional form  $f(x, \mathbf{w})$ . Ideally we want to perform inference directly in function space.*

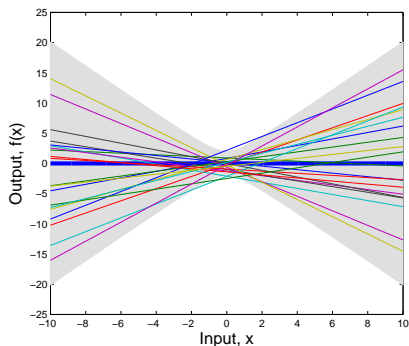
# The Function-Space View

A distribution over parameters  $p(\mathbf{w})$  induces a *distribution over functions*  $p(f(x))$ .

Consider the simple linear model,

$$f(x, w) = w_0 + w_1 x, \quad (1)$$

$$w_0, w_1 \sim \mathcal{N}(0, 1). \quad (2)$$



# The Function-Space View

We can express the distribution over functions *directly*. Suppose:

$$f(x, \mathbf{w}) = \mathbf{w}^T \phi(x) \quad (3)$$

$$p(\mathbf{w}) = \mathcal{N}(0, \Sigma_w) \quad (4)$$

$\phi(x)$  is a vector of basis functions. For example, if  $\phi(x) = (1, x, x^2)$  and  $x \in \mathbb{R}^1$  then  $f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2$  is a quadratic function.

**Can you derive the moments of the induced distribution over functions?**

$$\mathbb{E}[f(x, \mathbf{w})] = m(x) = \mathbb{E}[\mathbf{w}^T] \phi(x) =? \quad (5)$$

$$\text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = \mathbb{E}[f(x_i)f(x_j)] - \mathbb{E}[f(x_i)]\mathbb{E}[f(x_j)] =? \quad (6)$$

# The Function-Space View

We can express the distribution over functions *directly*. Suppose:

$$f(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x) \quad (7)$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_w) \quad (8)$$

## Moments of Induced Distribution over Functions

$$\mathbb{E}[f(x, \mathbf{w})] = m(x) = \mathbb{E}[\mathbf{w}^T] \boldsymbol{\phi}(x) = 0 \quad (9)$$

$$\text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = \mathbb{E}[f(x_i)f(x_j)] - \mathbb{E}[f(x_i)]\mathbb{E}[f(x_j)] \quad (10)$$

$$= \boldsymbol{\phi}(x_i)^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \boldsymbol{\phi}(x_j) - 0 \quad (11)$$

$$= \boldsymbol{\phi}(x_i)^T \Sigma_w \boldsymbol{\phi}(x_j) \quad (12)$$

**Are there any higher moments?**

# The Function-Space View

We can express the distribution over functions *directly*. Suppose:

$$f(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x) \quad (13)$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_w) \quad (14)$$

## Moments of Induced Distribution over Functions

$$\mathbb{E}[f(x, \mathbf{w})] = m(x) = \mathbb{E}[\mathbf{w}^T] \boldsymbol{\phi}(x) = 0 \quad (15)$$

$$\text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = \mathbb{E}[f(x_i)f(x_j)] - \mathbb{E}[f(x_i)]\mathbb{E}[f(x_j)] \quad (16)$$

$$= \boldsymbol{\phi}(x_i)^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \boldsymbol{\phi}(x_j) - 0 \quad (17)$$

$$= \boldsymbol{\phi}(x_i)^T \Sigma_w \boldsymbol{\phi}(x_j) \quad (18)$$

$f(x) \sim \mathcal{GP}(m, k)$  is a *Gaussian process* with mean function  $m(x)$  and covariance function (aka kernel)  $k(x, x')$ .

**We can do inference directly over  $f(x)$  instead of  $\mathbf{w}$ .**

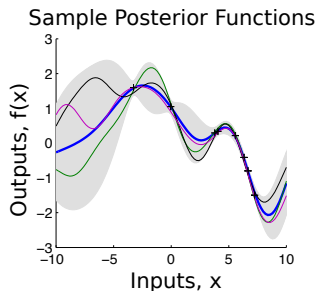
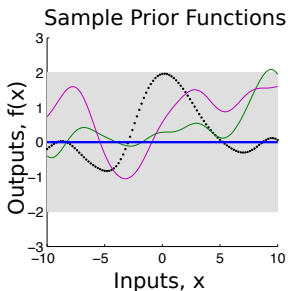
## Definition

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

## Nonparametric Regression Model

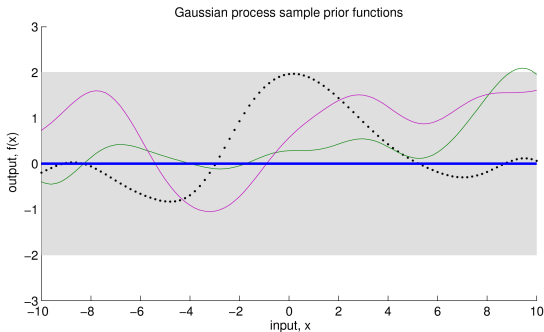
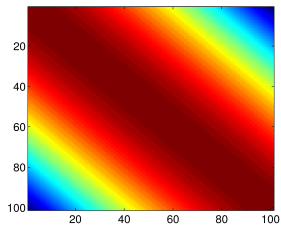
- Prior:  $f(x) \sim \mathcal{GP}(m(x), k(x, x'))$ , meaning  $(f(x_1), \dots, f(x_N)) \sim \mathcal{N}(\boldsymbol{\mu}, K)$ , with  $\boldsymbol{\mu}_i = m(x_i)$  and  $K_{ij} = \text{cov}(f(x_i), f(x_j)) = k(x_i, x_j)$ .

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



# Example: RBF Kernel

$$\begin{aligned}k_{\text{RBF}}(x, x') &= \text{cov}(f(x), f(x')) \\ &= a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)\end{aligned}$$



# Evaluating the RBF Kernel

Let  $\tau = x - x'$ :

$$k_{\text{RBF}}(\tau) = \exp(-0.5\tau^2/\ell^2) \quad (36)$$

Suppose we have input points  $x_1, x_2, x_3$ :

$K =$

$$\begin{pmatrix} \exp(-0.5(x_1 - x_1)^2/\ell^2) & \exp(-0.5(x_1 - x_2)^2/\ell^2) & \exp(-0.5(x_1 - x_3)^2/\ell^2) \\ \exp(-0.5(x_2 - x_1)^2/\ell^2) & \exp(-0.5(x_2 - x_2)^2/\ell^2) & \exp(-0.5(x_2 - x_3)^2/\ell^2) \\ \exp(-0.5(x_3 - x_1)^2/\ell^2) & \exp(-0.5(x_3 - x_2)^2/\ell^2) & \exp(-0.5(x_3 - x_3)^2/\ell^2) \end{pmatrix}$$

$=$

$$\begin{pmatrix} 1 & \exp(-0.5(x_1 - x_2)^2/\ell^2) & \exp(-0.5(x_1 - x_3)^2/\ell^2) \\ \exp(-0.5(x_2 - x_1)^2/\ell^2) & 1 & \exp(-0.5(x_2 - x_3)^2/\ell^2) \\ \exp(-0.5(x_3 - x_1)^2/\ell^2) & \exp(-0.5(x_3 - x_2)^2/\ell^2) & 1 \end{pmatrix}$$

$$K_{ij} = \text{cov}[f(x_i), f(x_j)].$$

# Sampling from GP Prior with RBF Kernel

```
x = [-10:0.2:10]'; % inputs (where we query the GP)
N = numel(x); % number of inputs
K = zeros(N,N); % covariance matrix

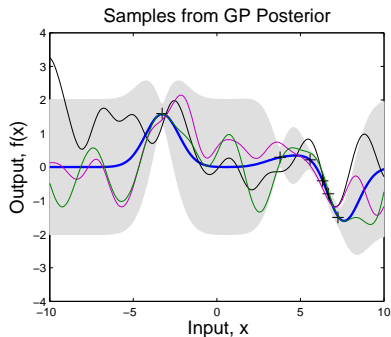
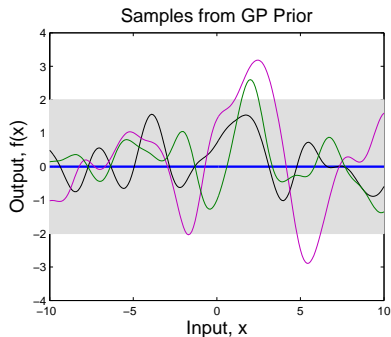
% very inefficient way of creating K in Matlab
for i=1:N
    for j=1:N
        K(i,j) = k_rbf(x(i),x(j));
    end
end

K = K + 1e-6*eye(N); % add jitter for conditioning of K
CK = chol(K);
f = CK'*randn(N,1); % draws from N(0,K)

plot(x,f);
```

# Inference using an RBF kernel

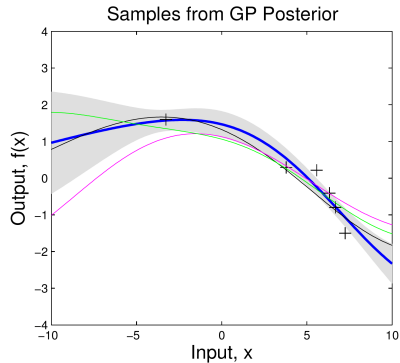
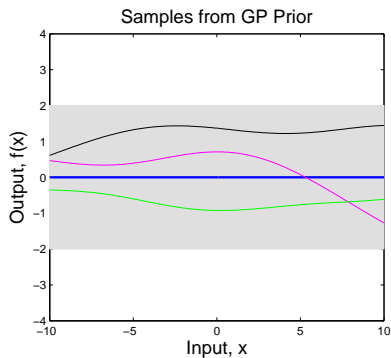
- ▶ Specify  $f(x) \sim \mathcal{GP}(0, k)$ .
- ▶ Choose  $k_{\text{RBF}}(x, x') = a_0^2 \exp(-\frac{\|x-x'\|^2}{2\ell_0^2})$ . Choose values for  $a_0$  and  $\ell_0$ .
- ▶ Observe data, look at the prior and posterior over functions.



- ▶ Does something look strange about these functions?

# Inference using an RBF kernel

Increase the length-scale  $\ell$ .



► Does something look strange about these functions?

# Inferring the Posterior

- ▶ 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

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

## Conditional predictive distribution

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

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

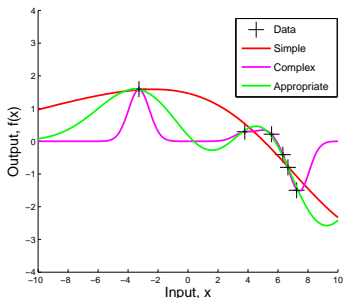
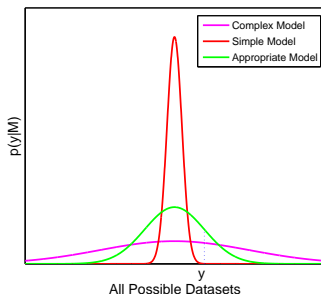
$$\text{cov}(\mathbf{f}_*) = K_\theta(X_*, X_*) - K_\theta(X_*, X) [K_\theta(X, X) + \sigma^2 I]^{-1} K_\theta(X, X_*).$$

# Learning and Model Selection

$$p(\mathcal{M}_i|\mathbf{y}) = \frac{p(\mathbf{y}|\mathcal{M}_i)p(\mathcal{M}_i)}{p(\mathbf{y})} \quad (19)$$

We can write the *evidence* of the model as

$$p(\mathbf{y}|\mathcal{M}_i) = \int p(\mathbf{y}|\mathbf{f}, \mathcal{M}_i)p(\mathbf{f})d\mathbf{f} \quad (20)$$



*Gaussian processes for Machine Learning*. Rasmussen, C.E. and Williams, C.K.I. MIT Press, 2006.

*Bayesian Methods for Adaptive Models*. MacKay, D.J. PhD Thesis, 1992.

*Covariance Kernels for Fast Automatic Pattern Discovery and Extrapolation with Gaussian Processes*. Wilson, A.G. PhD Thesis, 2014.

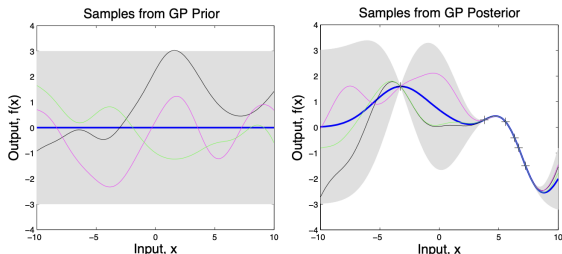
# Kernel Learning

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

$$p(\mathbf{y}|\boldsymbol{\theta}, X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|\boldsymbol{\theta}, X)d\mathbf{f}. \quad (32)$$

$$\log p(\mathbf{y}|\boldsymbol{\theta}, X) = \underbrace{-\frac{1}{2}\mathbf{y}^T(K_{\boldsymbol{\theta}} + \sigma^2 I)^{-1}\mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2}\log |K_{\boldsymbol{\theta}} + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2}\log(2\pi). \quad (33)$$

- ▶ An extremely powerful mechanism for kernel learning.



# Full Procedure

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

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

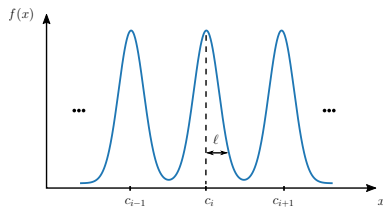
with respect to kernel hyperparameters  $\boldsymbol{\theta}$ . *The marginal likelihood provides a powerful mechanism for kernel learning.*

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

$$\begin{aligned} \mathbf{f}_*|X_*, X, \mathbf{y}, \boldsymbol{\theta} &\sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \\ \bar{\mathbf{f}}_* &= K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}\mathbf{y}, \\ \text{cov}(\mathbf{f}_*) &= K_{\boldsymbol{\theta}}(X_*, X_*) - K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}K_{\boldsymbol{\theta}}(X, X_*). \end{aligned}$$

$(K_{\boldsymbol{\theta}} + \sigma^2 I)^{-1}\mathbf{y}$  and  $\log |K_{\boldsymbol{\theta}} + \sigma^2 I|$  **naively require  $\mathcal{O}(n^3)$  computations,  $\mathcal{O}(n^2)$  storage.**

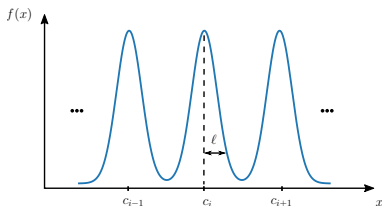
# Deriving the RBF Kernel



$$f(x) = \sum_{i=1}^J w_i \phi_i(x), \quad w_i \sim \mathcal{N}\left(0, \frac{\sigma^2}{J}\right), \quad \phi_i(x) = \exp\left(-\frac{(x - c_i)^2}{2\ell^2}\right) \quad (21)$$

$$\therefore k(x, x') = \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') \quad (22)$$

# Deriving the RBF Kernel



$$f(x) = \sum_{i=1}^J w_i \phi_i(x), \quad w_i \sim \mathcal{N}\left(0, \frac{\sigma^2}{J}\right), \quad \phi_i(x) = \exp\left(-\frac{(x - c_i)^2}{2\ell^2}\right) \quad (23)$$

$$\therefore k(x, x') = \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') \quad (24)$$

- Let  $c_J = \log J$ ,  $c_1 = -\log J$ , and  $c_{i+1} - c_i = \Delta c = 2\frac{\log J}{J}$ , and  $J \rightarrow \infty$ , the kernel in Eq. (27) becomes a Riemann sum:

$$k(x, x') = \lim_{J \rightarrow \infty} \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') = \int_{c_0}^{c_\infty} \phi_c(x) \phi_c(x') dc \quad (25)$$

# Deriving the RBF Kernel

$$f(x) = \sum_{i=1}^J w_i \phi_i(x), \quad w_i \sim \mathcal{N}\left(0, \frac{\sigma^2}{J}\right), \quad \phi_i(x) = \exp\left(-\frac{(x - c_i)^2}{2\ell^2}\right) \quad (26)$$

$$\therefore k(x, x') = \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') \quad (27)$$

- ▶ Let  $c_J = \log J$ ,  $c_1 = -\log J$ , and  $c_{i+1} - c_i = \Delta c = 2\frac{\log J}{J}$ , and  $J \rightarrow \infty$ , the kernel in Eq. (27) becomes a Riemann sum:

$$k(x, x') = \lim_{J \rightarrow \infty} \frac{\sigma^2}{J} \sum_{i=1}^J \phi_i(x) \phi_i(x') = \int_{c_0}^{c_\infty} \phi_c(x) \phi_c(x') dc \quad (28)$$

- ▶ 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:

$$k(x, x') = \int_{-\infty}^{\infty} \exp\left(-\frac{(x - c)^2}{2\ell^2}\right) \exp\left(-\frac{(x' - c)^2}{2\ell^2}\right) dc \quad (29)$$

$$= \sqrt{\pi} \ell \sigma^2 \exp\left(-\frac{(x - x')^2}{2(\sqrt{2}\ell)^2}\right) \propto k_{\text{RBF}}(x, x'). \quad (30)$$

# 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.
- ▶ GPs with RBF kernels are **flexible** (universal approximators) but **simple**.
- ▶ These models have strong *inductive biases* that concentrate prior support around simple solutions, providing good generalization even on small datasets, while retaining flexibility.

# 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:

$$f(x) = b + \sum_{i=1}^J v_i h(x; \mathbf{u}_i). \quad (31)$$

- ▶  $b$  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  $\mathbf{w}$ ,

$$\mathbb{E}_{\mathbf{w}}[f(x)] = 0, \quad (32)$$

$$\text{cov}[f(x), f(x')] = \mathbb{E}_{\mathbf{w}}[f(x)f(x')] = \sigma_b^2 + \frac{1}{J} \sum_{i=1}^J \sigma_v^2 \mathbb{E}_{\mathbf{u}}[h_i(x; \mathbf{u}_i)h_i(x'; \mathbf{u}_i)], \quad (33)$$

$$= \sigma_b^2 + \sigma_v^2 \mathbb{E}_{\mathbf{u}}[h(x; \mathbf{u})h(x'; \mathbf{u})]. \quad (34)$$

We can show any collection of values  $f(x_1), \dots, f(x_N)$  must have a joint Gaussian distribution using the central limit theorem.

# Neural Network Kernel

$$f(x) = b + \sum_{i=1}^J v_i h(x; \mathbf{u}_i). \quad (35)$$

- ▶ 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

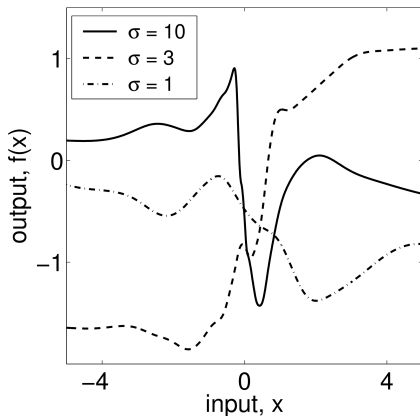
$$k_{\text{NN}}(x, x') = \frac{2}{\pi} \sin\left(\frac{2\tilde{x}^T \Sigma \tilde{x}'}{\sqrt{(1 + 2\tilde{x}^T \Sigma \tilde{x})(1 + 2\tilde{x}'^T \Sigma \tilde{x}')}}\right), \quad (36)$$

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

# Neural Network Kernel

$$k_{\text{NN}}(x, x') = \frac{2}{\pi} \sin\left(\frac{2\tilde{x}^T \Sigma \tilde{x}'}{\sqrt{(1 + 2\tilde{x}^T \Sigma \tilde{x})(1 + 2\tilde{x}'^T \Sigma \tilde{x}')}}\right) \quad (37)$$

Set  $\Sigma = \text{diag}(\sigma_0, \sigma)$ . Draws from a GP with a neural network kernel with varying  $\sigma$ :

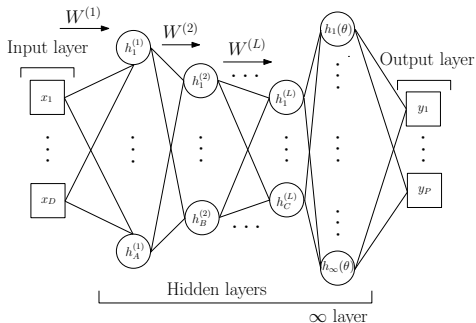


“How can Gaussian processes possibly replace neural networks? Have we thrown the baby out with the bathwater?” (MacKay, 1998)

*Introduction to Gaussian processes.* MacKay, D. J. In Bishop, C. M. (ed.), *Neural Networks and Machine Learning*, Chapter 11, pp. 133-165. Springer-Verlag, 1998.

# Deep Kernel Learning

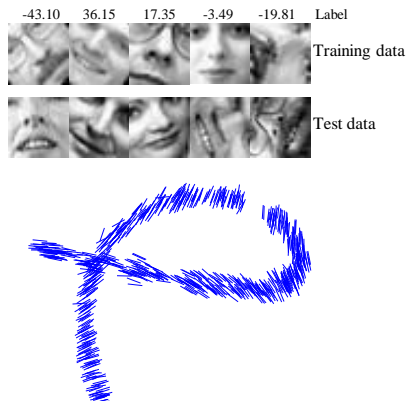
*Deep kernel learning* combines the inductive biases of deep learning architectures with the non-parametric flexibility of Gaussian processes.



**Base kernel hyperparameters  $\theta$  and deep network hyperparameters  $w$  are jointly trained through the marginal likelihood objective.**

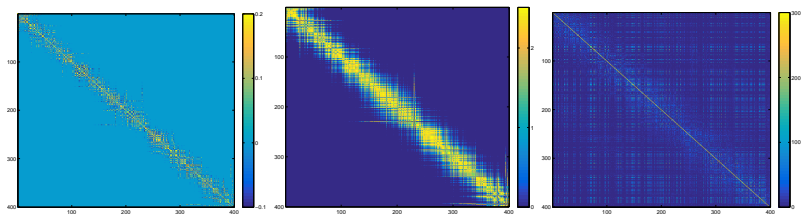
*Deep Kernel Learning*. Wilson, A.G., Hu, Z., Salakhutdinov, R., Xing, E.P. AISTATS, 2016

# Face Orientation Extraction



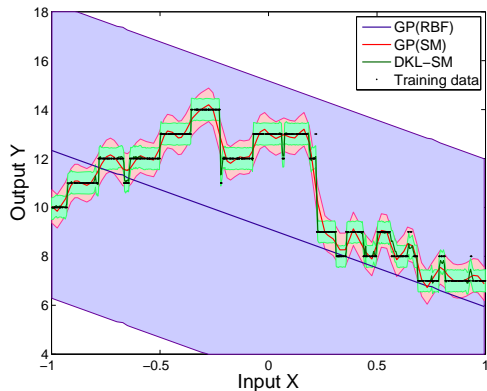
**Figure:** **Top:** Randomly sampled examples of the training and test data. **Bottom:** The two dimensional outputs of the convolutional network on a set of test cases. Each point is shown using a line segment that has the same orientation as the input face.

# Learning Flexible Non-Euclidean Similarity Metrics



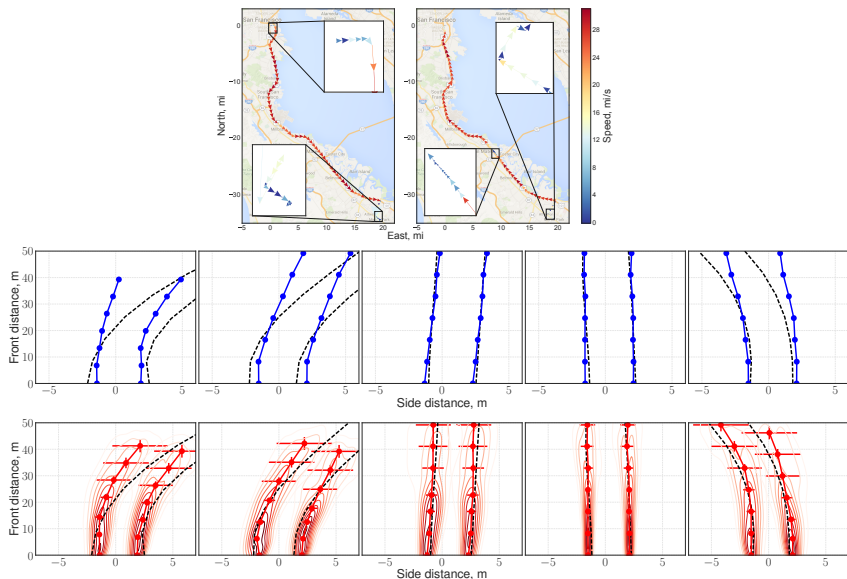
**Figure:** **Left:** The induced covariance matrix using DKL-SM (spectral mixture) kernel on a set of test cases, where the test samples are ordered according to the *orientations* of the input faces. **Middle:** The respective covariance matrix using DKL-RBF kernel. **Right:** The respective covariance matrix using regular RBF kernel. The models are trained with  $n = 12,000$ .

# Step Function



**Figure:** Recovering a step function. We show the predictive mean and 95% of the predictive probability mass for regular GPs with RBF and SM kernels, and DKL with SM base kernel.

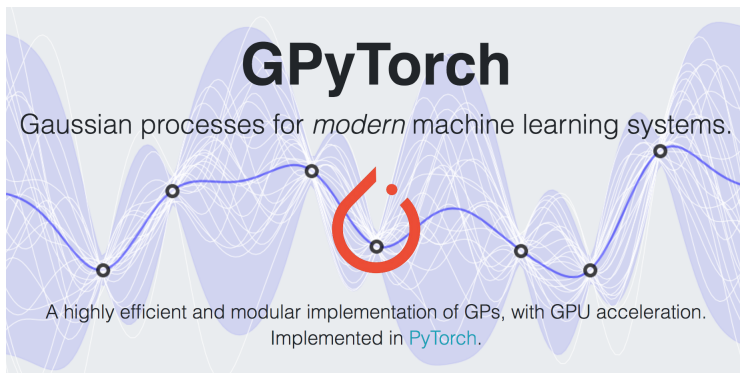
# Deep Kernel Learning for Autonomous Driving



*Learning scalable deep kernels with recurrent structure.*  
Al-Shedivat, M., Wilson, A.G., Saatchi, Y., Hu, Z., Xing, E.P. JMLR, 2017.

# Scalable Gaussian Processes

- ▶ Run **exact** GPs on millions of points in minutes.
- ▶ Based on Krylov subspace methods that permit significant GPU acceleration.
- ▶ Outperforms stand-alone deep neural networks by learning **deep kernels**.
- ▶ **Implemented in the new library GPyTorch:** `gpytorch.ai`

The image shows the GPyTorch logo, which consists of a stylized red 'G' with a dot above it, set against a background of blue wavy lines representing Gaussian process uncertainty. The text 'GPyTorch' is prominently displayed in large black font, with the tagline 'Gaussian processes for modern machine learning systems.' below it. At the bottom, it states 'A highly efficient and modular implementation of GPs, with GPU acceleration. Implemented in PyTorch.'

**GPyTorch**  
Gaussian processes for *modern* machine learning systems.

A highly efficient and modular implementation of GPs, with GPU acceleration.  
Implemented in [PyTorch](#).

*GPyTorch: Blackbox Matrix-Matrix Gaussian Process Inference with GPU Acceleration.*  
Gardner, J. R., Pleiss, G., Bindel, D., Weinberger, K. Q., Wilson, A. G. NeurIPS, 2018.