

## CUED - Engineering Tripos Part IIB 2024-2025

## Module Coursework

Module	4F13	Title of report	Gaussian Processes
Date submitted: 6/11/24		Assessment for this module is <input checked="" type="checkbox"/> 100% / <input type="checkbox"/> 25% coursework of which this assignment forms <u>33.3</u> %	
<b>UNDERGRADUATE and POST GRADUATE STUDENTS</b>			
Candidate number:	5488A		<input checked="" type="checkbox"/> Undergraduate <input type="checkbox"/> Post graduate

## Feedback to the student

☐ See also comments in the text

Feedback to the student		Very good	Good	Needs improvmt
C O N T E N T	<b>Completeness, quantity of content:</b> Has the report covered all aspects of the lab? Has the analysis been carried out thoroughly?			
	<b>Correctness, quality of content</b> Is the data correct? Is the analysis of the data correct? Are the conclusions correct?			
	<b>Depth of understanding, quality of discussion</b> Does the report show a good technical understanding? Have all the relevant conclusions been drawn?			
	Comments:			
P R E S E N T A T I O N	<b>Attention to detail, typesetting and typographical errors</b> Is the report free of typographical errors? Are the figures/tables/references presented professionally?			
	Comments:			

Marker:

Date:

# 4F13 Coursework 1 - Gaussian Processes

November 2024

## 1 Task A

Our model is defined in Equation 1, we have a gaussian process ( $f$ ) with a squared exponential (SE) covariance function , Equation 2, zero mean and gaussian likelihood.

$$y = f(x) + \eta : f \sim \mathcal{N}(0, k_{SE}(x, x')); \eta \sim \mathcal{N}(0, \sigma_n^2) \quad (1)$$

$$k_{SE}(x, x') = \sigma_f^2 \exp\left(-\frac{(x - x')^2}{2\lambda^2}\right) \quad (2)$$

The model hyper-parameters are trained by minimising the negative log marginal likelihood ( $\mathcal{L}$ ). We do this and generate the predictive distribution using the code in Listing 1.

```
meanf = []; covf = @covSEiso; likf = @likGauss;
hyp_init.mean = []; hyp_init.cov = [-1 0]; hyp_init.lik = 0;
hyp_opt = minimize(hyp_init, @gp, -100, @infGaussLik, meanf, covf, likf, x, y);
[mu, s2] = gp(hyp_opt, @infGaussLik, meanf, covf, likf, x, y, xs);
```

Listing 1: Code to train hyper-parameters and generate the predictive distribution of a GP with squared exponential covariance

The trained hyper-parameters are listed as Optimum 1 in Table 1. We plot the data and predictive distribution in Figure 1a. The hyper-parameters have the following interpretation;  $\lambda$  - length scale,  $\sigma_f$  - scale factor and  $\sigma_n^2$  - measurement noise variance.

	$\lambda$	$\sigma_f$	$\sigma_n$	$\mathcal{L}$
Optimum 1	0.128	0.897	0.118	$1.19 \times 10^1$
Optimum 2	8.049	0.696	0.663	$7.82 \times 10^1$

Table 1: Hyper-parameter values at 2 local minima of  $\mathcal{L}$

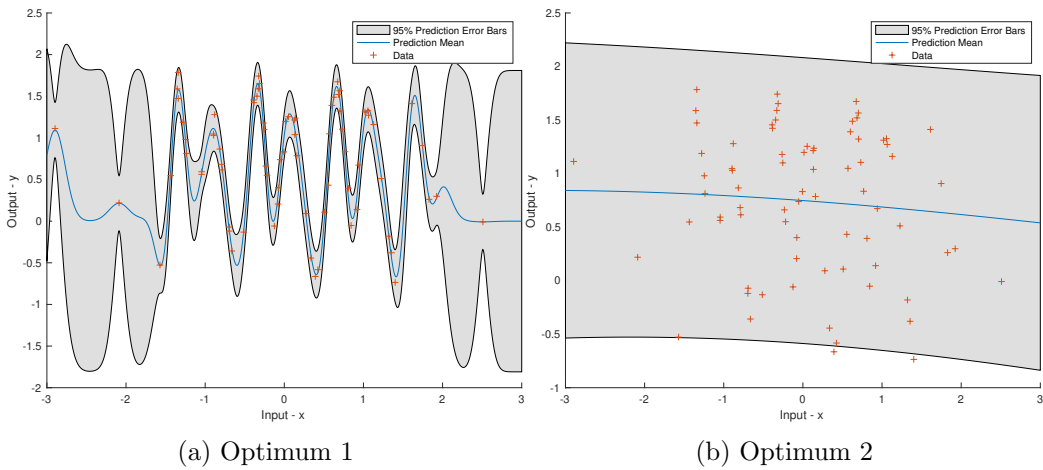


Figure 1: Predictive mean and 95% error bars of models with hyper-parameters given in Table 1. Training data is plotted as well.

From Figure 1a we see that in regions of higher data density the error bars are small while where data is sparse they become wider and approach a constant value. This is

explained by the form of the predictive variance, Equation 3. Our data inputs are denoted  $\mathbf{x}$ . (When calculating the predictive error bars we evaluate our predictive covariance with  $x' = x$ , Equation 3 gives the simplified equation in this case.)

$$k_{|y}(x) = \sigma_f^2 + \sigma_n^2 - k_{SE}(x, \mathbf{x})[k_{SE}(\mathbf{x}, \mathbf{x}) + \sigma_n^2 I]^{-1} k_{SE}(\mathbf{x}, x) \quad (3)$$

There are 3 terms, the first two are constant,  $\sigma_f^2 + \sigma_n^2$ , these are the prior variance.  $k_{SE}(\mathbf{x}, \mathbf{x}) + \sigma_n^2 I$  is positive definite by definition of the covariance kernel therefore the third term is always negative. This third term becomes larger in magnitude when there are many data points within  $\lambda$  of  $x$ , this decreases the predictive variance where there is a higher density of data. This form makes intuitive sense too, we can be more confident in predictions where we have more data and where we have no data we can only use our prior knowledge.

## 2 Task B

To identify all local minima of  $\mathcal{L}$  we perform a grid search over the hyper-parameters. Figure 2 shows a contour  $\mathcal{L}$  for a slice of the search with  $\sigma_f = 1$  which shows the two minima well. The second optimum (Optimum 2 in Table 1) has a much longer length scale,  $\lambda$ , and larger measurement noise,  $\sigma_n$ . From the predictive distribution, Figure 1b, we see that this optimum results in a model that explains most of the output variation as measurement noise instead of the value of the function unlike Optimum 1 which does the opposite. However, this second optimum is a worse fit with a higher value of  $\mathcal{L}$ . Furthermore, by observing the distribution of the residuals we see that they do not seem to be independent of the input variable, while our model expects independent measurement noise. Therefore, we can conclude that Optimum 1 is more likely to be the model that generated this data.

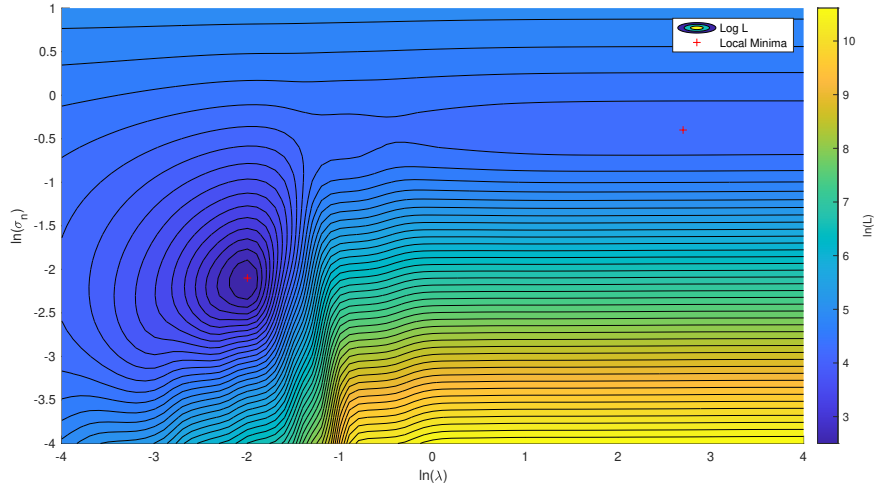


Figure 2: Contour of  $\ln(\mathcal{L})$  with  $\lambda$  and  $\sigma_n$  for fixed  $\sigma_f = 1$ . Shows the two local minima of  $\mathcal{L}$  given in Table 1. We are plotting  $\ln(\mathcal{L})$  for more evenly spaced contours.

## 3 Task C

```
covf = @covPeriodic; hyp_init.cov = [-1 0 0];
```

Listing 2: Code to use periodic SE covariance. Training and prediction code same as Listing 1

Equation 4 gives the form of the periodic squared exponential covariance function. It has a very similar form to the standard squared exponential covariance function, but the

measure of "distance" between two points in input space is now  $\sin(\frac{\pi}{p}(x - x'))$ . This means that the "distance" between two points any multiple of  $p$  apart is now zero, giving large covariance between these points. This means that samples from this GP will be periodic with period  $p$ . The optimised hyper-parameters for the periodic SE covariance function are given in Table 2 and the prediction intervals are shown in Figure 3.

$$k_{PSE}(x, x') = \sigma_f^2 \exp(-\frac{2}{\lambda^2} \sin^2(\frac{\pi}{p}(x - x')))) \quad (4)$$

$\lambda$	$p$	$\sigma_f$	$\sigma_n$	$\mathcal{L}$
0.705	0.999	0.694	0.085	$-2.93 \times 10^1$

Table 2: Hyper-parameter values for periodic SE covariance function

The effect of the periodic covariance is clear in the prediction intervals. Unlike the previous model where prediction intervals were large in where the density of data was lower, now, as long as there is data a multiple of the period apart, the model has small prediction intervals. This is due to the form of the covariance function and predictive distribution.

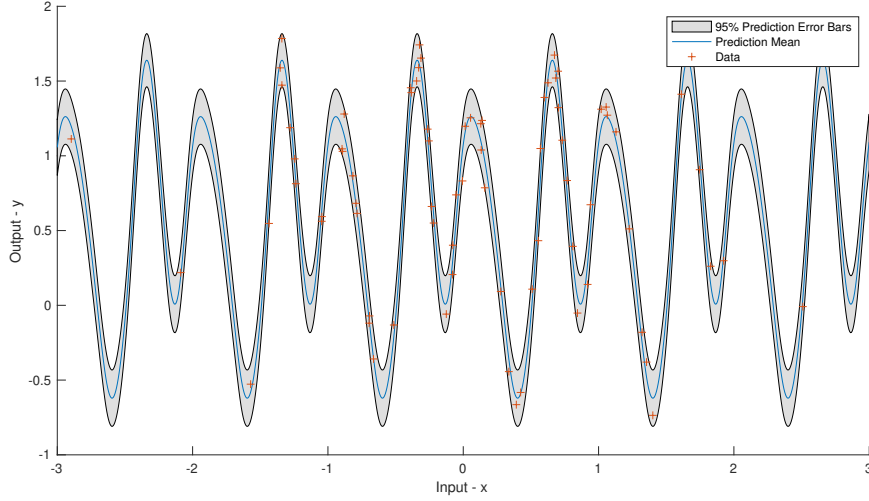


Figure 3: Prediction intervals for periodic SE covariance function with hyper-parameters given in Table 2

The marginal likelihood of the periodic model is higher than that of the standard squared exponential model, indicating a better fit. We also gain confidence in this model as our training data includes multiple periods. Further evidence that the periodic model is accurate can be seen in the residuals of the data, which are close normally distributed, Figure 4a, and independent of the input variable, Figure 4b. This matches our model definition.

## 4 Task D

```
N_points = 200; N_samples = 2;
x = linspace(-5,5,N_points)';
K = feval(covf{:}, hyp.cov, x);
y = chol(K + 1e-6*eye(N_points))' * gpml_randn(2, N_points, N_samples);
```

Listing 3: Code to generate samples from a GP with covariance given by covf

To sample from a GP we can choose a set of evaluation points -  $\mathbf{x}$  and evaluate our mean and covariance functions at these points to obtain a mean vector -  $\boldsymbol{\mu}$  and covariance matrix -  $K$ . We can generate our samples  $\mathbf{y} = \boldsymbol{\mu} + \text{chol}(K)^T \boldsymbol{\eta}$  where  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, I)$ .

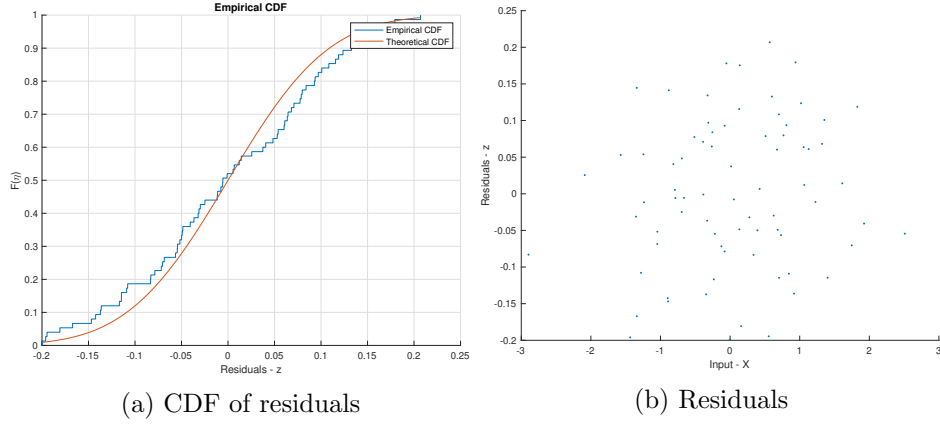


Figure 4: Residuals of data when fitted with periodic SE covariance GP

We require the Cholesky decomposition of  $K - \text{chol}(K)$ , however, this is only possible for positive definite matrices. Our covariance kernel is positive semi-definite, therefore, to ensure  $K$  is positive definite we can add  $\epsilon I$  before the decomposition. This ensures it is positive definite without changing the behaviour of the matrix much. The code to generate the samples is shown in Listing 3 for  $\epsilon = 1 \times 10^{-6}$ .

We sample a GP with a covariance kernel that is the product of a long length scale SE kernel and short length scale periodic SE kernel. The form of the kernel is given in Equation 5 and hyper-parameter values are given in Table 3.

$$k(x, x') = k_{PSE|p^{(1)}, \lambda^{(1)}, \sigma_f^{(1)}}(x, x') k_{SE|\lambda^{(2)}, \sigma_f^{(2)}}(x, x') \quad (5)$$

$p^{(1)}$	$\lambda^{(1)}$	$\sigma_f^{(1)}$	$\lambda^{(2)}$	$\sigma_f^{(2)}$
1	0.607	1	7.389	1

Table 3: Hyper-parameter values for periodic SE covariance function

We observe that the samples in Figure 5c exhibit characteristics of both kernels shown in Figures 5a and 5b. Specifically, they display periodic behaviour with a period of  $p^{(1)} = 1$  over short length scales, like Figure 5a, combined with gradual changes over larger length scales, Figure 5b.

In comparison to a purely periodic model, the prediction intervals of this model will not remain narrow even very far from any data. They will instead slowly decay with the length scale of the SE component. This strikes a middle ground between the models investigated in Task A and C where we want to gain prediction accuracy using the periodic behaviour while not making narrow predictions very far outside of our training data.

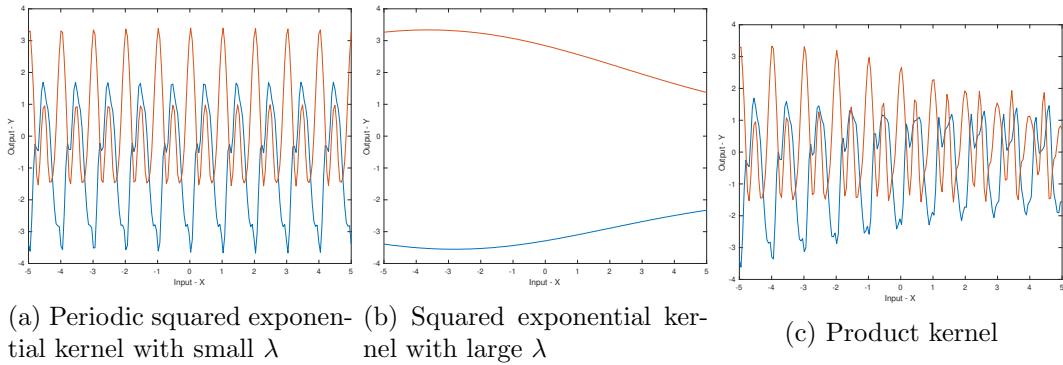


Figure 5: Samples from 3 covariance kernels, hyper-parameters for each given in Table 3

## 5 Task E

The squared exponential Automatic Relevance Determination (SE-ARD) kernel, given in Equation 6, is similar to the standard squared exponential kernel, but the distance measure can now be weighted differently for each input dimension. This can be useful when the input dimensions have different units or scales.

$$k_{SE-ARD}(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{i=1}^D \frac{(x_i - x'_i)^2}{\lambda_i}\right) \quad (6)$$

When we fit the sample data with the SE-ARD kernel we get the hyper-parameters given in Table 4 case A. The code to do this is similar to Listing 1. The length scales in each input dimension are similar so the model is not making use of the ARD property. The prediction intervals of the model are shown in Figure 6a.

Case	$\sigma_f^1$	$\lambda_1^1$	$\lambda_2^1$	$\sigma_f^1$	$\lambda_1^1$	$\lambda_2^1$	$\sigma_n$	$\mathcal{L}$
A: $k_{SE-ARD}$	1.107	1.511	1.285	-	-	-	0.1026	$-1.9218 \times 10^1$
B: $k_{SE-ARD}^{(1)} + k_{SE-ARD}^{(2)}$	0.7116	1104	0.9864	1.108	1.446	1281	0.0979	$-6.6394 \times 10^1$

Table 4: Hyper-parameter values for periodic SE covariance function

When our kernel is the sum of two independent SE-ARD kernels case B we observe that the fitted hyper parameters change dramatically. In each SE-ARD kernel, one length scale parameter is significantly larger than the other. Therefore, that dimension has almost no effect on the covariance between points. This is similar to a covariance function that is the sum of two scalar squared exponential kernels in each input dimension. Essentially our GP can be decomposed as  $f(x, y) = f_x(x) + f_y(y)$ .

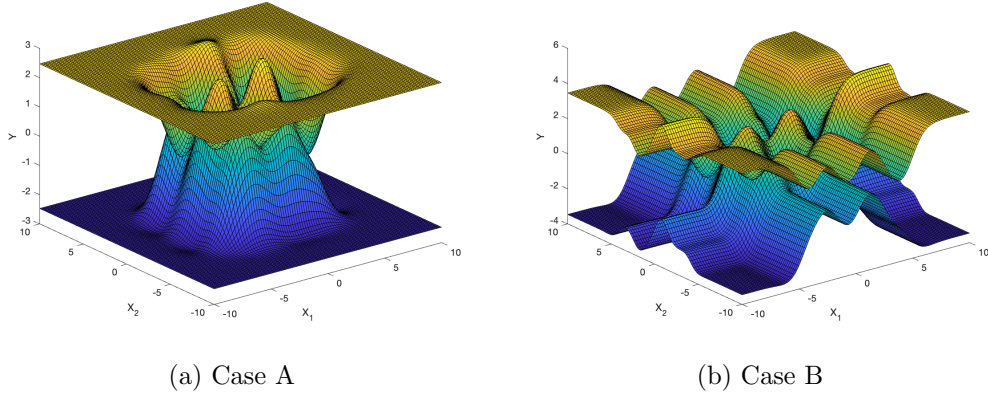


Figure 6: 95% Prediction intervals for models given in Table 4 using data from cw1e.mat

Both models have similar prediction intervals in the region of the data. However, the behaviour of the prediction intervals when extrapolating is different. While the case A very quickly returns to the prior prediction interval, the case B keeps smaller prediction intervals in directions parallel to the input axes. This is because in these directions, the "distance" measure to the data in one of the two SE-ARD kernels is small, as the corresponding ARD  $\lambda$  is very large. This results in a decrease in the prediction interval from that kernel.

The model in case B has a higher likelihood than the one in case A however we would argue that the case A model is better in most cases. By observing our dataset, we have not covered enough of the input space to conclude that we can decompose the function into a sum of functions in each input dimension. It would therefore not be appropriate to be confident in the extrapolation behaviour of this model.