

FYP

Sean White

January 2025

1 GR intro

Gravitational waves (GWs) are small fluctuations of spacetime that propagate at the speed of light. As described in *Gravitational Waves, Vol. 1* by Maggiore [4, p. 29],

“In this setting, the definition of GWs is relatively clear: the background spacetime is flat, and the small fluctuations around it have been called ‘gravitational waves’. The term ‘waves’ is justified by the fact that, in a suitable gauge, $h_{\mu\nu}$ indeed satisfies a wave equation.”

This is formally approached by expanding the Einstein equations around the flat Minkowski metric η_{ab}

$$g_{ab} = \eta_{ab} + \epsilon h_{ab}, \quad \epsilon \ll 1, \quad \eta_{ab} = \text{diag}(-1, 1, 1, 1), \quad (1)$$

where h_{ab} represents the perturbation. In the linearized regime of general relativity, the Einstein field equations simplify to

$$\square \bar{h}_{ab} = \frac{-16\pi G}{c^4} T_{\mu\nu}. \quad (2)$$

However we are interested in this equation outside of the source (i.e $T_{\mu\nu} = 0$) therefore we are left with

$$\square \bar{h}_{ab} = 0, \quad (3)$$

with \square the d'Alembertian operator in flat spacetime.

Although h_{ab} initially has 10 independent components (as a symmetric 4×4 tensor), 8 of these correspond to gauge freedom and constraints. After imposing the Lorentz and transverse-traceless (TT) gauge conditions, only two physical degrees of freedom remain: the plus (h_+) and cross (h_\times) polarizations [4, Sec. 1.2].

This wave equation admits plane-wave solutions. For a wave propagating in the z -direction, the TT-gauge form of the perturbation is:

$$h_{ab}^{(\text{TT})} \propto \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & h_+ & h_\times & 0 \\ 0 & h_\times & -h_+ & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} e^{i(kz - \omega t)}. \quad (4)$$

Leading-Order Power Emission by Gravitational Waves

In general, the power emitted by a radiative field can be expressed schematically as

$$\dot{E} = \sum_{\ell=0}^{\infty} \left\langle \left| \left(\frac{\partial}{\partial t} \right)^{\ell+1} P_\ell(t) \right| \right\rangle. \quad (5)$$

Here, $P_\ell(t)$ represents the multipole moments of the radiating source. The $\ell = 0$ term corresponds to the monopole moment, which in the gravitational case is the total mass of the system. Assuming mass is conserved, this term vanishes.

The $\ell = 1$ term represents the dipole moment, which is also zero in the gravitational case due to conservation of linear momentum. Therefore, the leading-order contribution to gravitational-wave emission arises from the $\ell = 2$ term, known as the quadrupole.

This gives the leading-order expression for the power emitted in gravitational waves:

$$\dot{E} = \frac{G}{5c^5} \left\langle \ddot{Q}_{ij} \ddot{Q}^{ij} \right\rangle, \quad (6)$$

where Q_{ij} is the mass quadrupole moment of the source. It is related to the mass moment M_{ij} by

$$Q_{ij} := M_{ij} - \frac{1}{3}\delta_{ij}M_k^k, \quad (7)$$

where $M_{ij} = \int d^3x T^{00}(x^i x^j)$.

A Simple Two-Mass Inspiral Model

I took this example and summarised it from Sarp's notes, leads nicely to mismatch

We will now consider a generic problem as illustrated in [1] where two point masses, $m_1 \geq m_2$, are in a quasi-circular orbit of separation R , each at distances r_1 and r_2 from their common center of mass (CoM), with $R = r_1 + r_2$. We place the orbit in the x - y plane so that mass 1 moves on $\mathbf{x}_1(t) = r_1(\cos \Omega t, \sin \Omega t)$ and mass 2 on $\mathbf{x}_2(t) = r_2(\cos(\Omega t + \pi), \sin(\Omega t + \pi))$. The system's orbital frequency is Ω . A short calculation yields

$$Q^{ij} = 4\Omega^3 (m_1 r_1^2 + m_2 r_2^2) \begin{pmatrix} \sin(2\Omega t) & -\cos(2\Omega t) \\ -\cos(2\Omega t) & -\sin(2\Omega t) \end{pmatrix}. \quad (8)$$

Introducing the reduced mass $\mu = m_1 m_2 / (m_1 + m_2)$, we get that the $\ell = 2$ power emission is

$$\dot{E}_{\ell=2} = \frac{32}{5} \frac{G}{c^5} \Omega^6 \mu^2 R^4. \quad (9)$$

Both Ω and R are functions of time, but they evolve on a time scale much longer than the orbital time scale and so when averaging over orbits we approximate them as constant. Applying Kepler's law, $\Omega^2 = GM/R^3$, and defining $\omega = 2\Omega$ as the GW frequency, we find

$$\dot{E}_{\ell=2} = \frac{32}{5} \frac{c^5}{G} \left(\frac{G \mathcal{M} \omega}{2c^3} \right)^{10/3}, \quad \text{where } \mathcal{M} = \mu^{3/5} M^{2/5} \quad (10)$$

is known as the chirp mass where $M = m_1 + m_2$ is the total mass. SA: You already kind of mention this below Eq. (9), so maybe provide slightly more detail there.

Could maybe mention radiation time scale being much smaller than orbit timescale but I don't see what this will add

\dot{E} represents the rate at which the system loses energy due to gravitational-wave emission. This energy loss causes the binary orbit to shrink and the GW frequency to increase, with the chirp mass \mathcal{M} and frequency ω capturing the key features of this inspiral.

SA: From here you should go on to derive the form of a simple waveform in first time domain then frequency domain. Then you can go into the Mismatch

Introducing the Waveform Mismatch

This simplified two-mass, circular-orbit model captures the main physical components of the model such as the "chirp behaviour" (frequency of GW increasing as the orbit shrinks). However we want to quantify how neglecting for example eccentricity of the orbit and the spin vectors of the masses effects the waveform. We compare the

simplified model to more complete waveforms using the waveform mismatch discussed in [5] and [3]. The mismatch between two signals is defined by

$$\mathcal{M} = 1 - \max_{\lambda_m} \frac{\langle h_{\text{simple}}, h_{\text{accurate}} \rangle}{\sqrt{\langle h_{\text{simple}}, h_{\text{simple}} \rangle \langle h_{\text{accurate}}, h_{\text{accurate}} \rangle}}, \quad (11)$$

where $\langle \cdot, \cdot \rangle$ is the inner product of both GW, and we maximise over a set of (intrinsic or extrinsic) model parameters λ_m . A mismatch $\mathcal{M} \ll 1$ indicates that our simple waveform faithfully represents the physical signal, whereas larger mismatches highlight missing physics (e.g. not circular orbit).

2 Background

This section was written ages ago and needs to be updated but I think could still be worth time just going over how bayesian methods are currently used in GW and how making a mismatch predictor model may help moving forward

The gold standard gravitational waveforms are obtained by numerical relativity simulations which involve directly solving Einstein's equations of general relativity which is massively computationally expensive. For this reason, only several thousand simulations are currently available. As a result, GW models rely on analytical or semi-analytical prescriptions that are calibrated to the numerical relativity simulations.

Each of these modelling approaches will incur a certain amount of error. We quantify this error as the mismatch between the "true" relativity simulation of a gravitational waveform and the analytic model for the waveform. The mismatch [5] varies between 0, signifying that the model and the true waveform are identical (up to an overall amplitude rescaling), and 1, meaning that the two are completely orthogonal. To begin we will delve into the methods used before to "combine" these models and then discuss the new proposed method. This will lead us onto my project.

Methods Used Before

Standard Method

We build a distribution based off of each of the 3 models. We then combine these distributions with equal weight to form a new total distribution from which we make our predictions from.

Evidence-Informed Method

This method again builds a distribution based off of each of the 3 models. The distributions are then combined with weights which are determined by how well the model explains the observed data. From this new distribution predictions are made.

New Proposed Method: Numerical Relativity Informed Method

Numerical Relativity Informed Method

This method builds a distribution based off of each of the 3 models. Then the mismatch (difference between NR simulations and model) is calculated. We then find the ratios of the mismatch between each model, for example:

$$\frac{\text{mismatch(model 1)}}{\text{mismatch(model 2)}}.$$

A distribution using the ratio of mismatch of each model is built. This distribution is used to determine the probability of selecting each model in each region of the parameter space. Predictions are then made using the model selected in that region of the parameter space using the ratio distribution. This selection of the model is probabilistic and so the best model in a certain region may be chosen 85% of the time and the second best 10% etc etc.

Comparison

The Standard Method and the Evidence-Informed Method build a distribution from which we make our predictions. The Numerical Relativity Informed method builds a distribution which helps us choose the best model in that parameter space “most of the time”.

Now the issue is that to calculate the mismatch for the models needed in the NR informed method, we must have gravitational waveforms generated from NR simulations. This as we discussed is a computationally expensive process and is often not feasible. The goal of my project is to build a Gaussian Process Regression model to model the mismatch of the model’s to the NR simulations.

3 Gaussian Process Regression GPR Background

3.1 Introduction and Roadmap

In the following subsections, I discuss the foundational concepts of Gaussian Process Regression (**GPR**) in four main steps:

1. **Gaussian Processes:** Section 3.2 introduces Gaussian Processes and explains how their priors and posteriors are constructed from finite sets of points.
2. **Kernel Functions:** In Section 3.3, I explore how kernels encode assumptions about smoothness, periodicity, and structural properties of the underlying function.
3. **Noise Modeling:** Section 3.4 covers several approaches for incorporating observational noise into the GP framework.
4. **Hyperparameter Optimization:** Finally, in Section 3.5, I discuss how kernel and noise hyperparameters influence the posterior and how they can be optimized to improve model performance.

3.2 Gaussian Proces Regression Background

Definition of a Gaussian Process

A Gaussian Process (**GP**) defines a probabilistic model over all possible functions rather than assuming a single function to be true

$$f(X) \sim \mathcal{GP}(\mu(X), k(X, X')). \quad (12)$$

where $\mu(X)$ is the mean function, specifying the expected function value at each X :

$$\mu(X) = E[f(X)], \quad (13)$$

$k(X, X')$ is the covariance function (kernel), encoding the relationships between function values at different points:

$$k(X, X') = \text{Cov}(f(X), f(X')). \quad (14)$$

Since the input space is continuous, the GP represents an infinite-dimensional distribution. In practice, we approximate the process by evaluating the GP at a finite set of inputs. These function values are then assumed to follow a multivariate normal Gaussian distribution.

Mathematically, for a finite set of input points

$$X = \{X_1, X_2, \dots, X_n\}, \quad (15)$$

the corresponding function values

$$f = \{f(X_1), f(X_2), \dots, f(X_n)\} \quad (16)$$

follow a multivariate normal distribution

$$f \sim \mathcal{N}(\mu(X), K(X, X)). \quad (17)$$

Each sample from this multivariate distribution represents a function evaluated at n different points.

The Prior Distribution

Before observing any data, we assume a joint Gaussian distribution over both training and test points. Let X denote training inputs and X_* test inputs. The joint prior over their function values is

$$\begin{bmatrix} f(X) \\ f(X_*) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu(X) \\ \mu(X_*) \end{bmatrix}, \underbrace{\begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}}_{C=\text{Covariance Matrix}} \right). \quad (18)$$

The corresponding joint probability density function (pdf) is given by:

$$p(f, f_*) = \frac{1}{(2\pi)^{n/2} \sqrt{|C|}} \exp \left(-\frac{1}{2} \left(\begin{bmatrix} f \\ f_* \end{bmatrix} - \begin{bmatrix} \mu(X) \\ \mu(X_*) \end{bmatrix} \right)^T C^{-1} \left(\begin{bmatrix} f \\ f_* \end{bmatrix} - \begin{bmatrix} \mu(X) \\ \mu(X_*) \end{bmatrix} \right) \right) \quad (19)$$

After accounting for the mean, the resulting distribution is entirely determined by its kernel function. The kernel governs how the model generalizes to unseen data. There are many kernel choices, each encoding different structural assumptions about the function, such as smoothness, periodicity, or linearity. In the next section we examine the different kernel choices available and the assumptions that each kernel encodes about our function structure, such as smoothness and periodicity.

3.3 Kernel Functions

The kernel function encodes our assumptions about the relationship between input points in a Gaussian Process (GP). It defines the covariance between any two function values and thereby determines the smoothness, periodicity, or other properties of the functions drawn from the GP prior. Fundamentally, kernels reflect the idea of *similarity*: input points x and x' that are close together are assumed to have highly correlated outputs $f(x)$ and $f(x')$, while distant inputs are assumed to produce less correlated values. This notion of similarity, as emphasized in [6, p. 79], is central to how Gaussian processes learn from and generalize beyond training data.

In Figure 1, we illustrate the effect of the kernel on the GP prior. We draw three functions from the multivariate Gaussian prior defined in Equation 18, using a zero mean and an RBF kernel. The first subplot shows these samples, while the second subplot visualizes the corresponding covariance matrix as a heatmap. The matrix reveals that correlations are strongest when input points are close together (near the diagonal) and decay as the distance between inputs increases. This is evident also from the samples as we can see nearby points often move in similar directions, while distant points diverge more significantly.

The final three subplots highlight how this distance-based correlation manifests in the joint distribution of pairs of function values. For closely spaced inputs, such as $(x, x') = (0, 0.1)$, the joint distribution of $(f(0), f(0.1))$ forms a narrow elliptical contour, indicating strong correlation (approximately 0.9). As the distance increases, such as in the pairs $(0, 0.5)$ and $(0, 1)$, the ellipses widen, reflecting weaker correlation. This visualization reinforces the intuition that kernel functions govern how input proximity translates to output similarity.

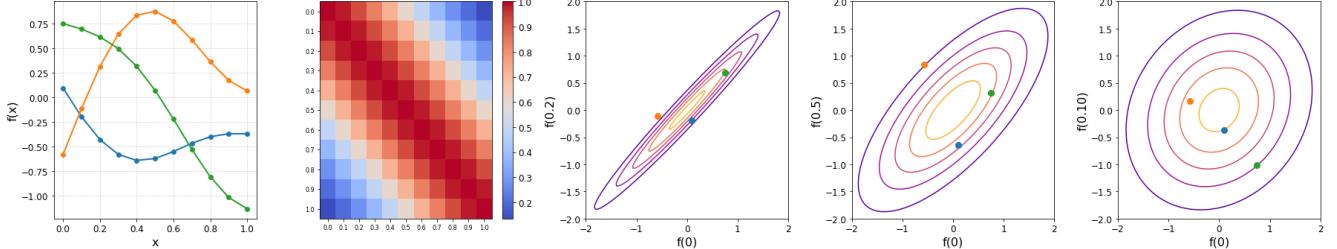


Figure 1: Sampling from the GP prior with zero mean and an RBF kernel ($\ell = 0.5$, $\sigma_f^2 = 1$). The first plot shows three sample functions drawn from the prior distribution. The second plot visualizes the covariance matrix as a heatmap, revealing the strength of correlations between inputs. The final three subplots display joint distributions between selected input pairs, illustrating how output correlation diminishes with increasing input distance.

We have discussed how the kernel function encodes the covariance structure of the GP prior. This structure depends on the choice of kernel. According to [6] kernels can be divided into two major sub-groups, stationary kernels and non-stationary kernels. Stationary kernels depend only on the relative (often radial) distance between inputs $\|x - x'\|$ and are invariant to translations in the input domain. By contrast, non-stationary kernels depend explicitly on the absolute values of x and x' , allowing the function’s properties—such as smoothness or amplitude—to vary across the domain. For more detailed discussion on building, combining, and customizing these kernels, see [2] and [6, Ch. 4].

In Table 1, we provide an overview of several common kernel types, showing both their functional form and samples drawn from the corresponding GP priors. While each kernel imposes a distinct structural pattern on the functions—such as smoothness, periodicity, or linearity—they are all similarly influenced by shared hyperparameters like the lengthscale. In addition, many kernels include unique internal parameters that further shape the behaviour of the modeled functions. In the following subsections, we explore each of these kernels in detail and discuss the role of their associated hyperparameters.

Kernel name:	RBF (SE)	Rational Quadratic	Periodic	Matern	Laplace	Linear (Dot Product)
Plot of $k(x, x')$:						
GP Prior Samples:						
Key Hyperparameters	ℓ (Lengthscale)	α (Scale-mix)	p (Period)	ν (Smoothness)	γ (Decay rate)	None or variance
Structure type:	Local variation	Multi-scale local variation	Repeating structure	Rough to smooth	Rougher variation	Linear functions

Table 1: Visual comparison of common kernel functions and their effect on Gaussian process priors. Each column shows the kernel shape $k(x, x')$, samples from the corresponding GP prior, and a summary of the structure it imposes. All kernels were evaluated using a lengthscale parameter $\ell = 1$ (except where noted). For the Matern kernel, $\nu = 0.5$; Laplace kernel, $\gamma = 6$; Rational Quadratic kernel, $\alpha = 0.25$; and Periodic kernel, period $p = 2$. Detailed formula and graphs for each kernel are provided in the appendix B.

Note: In practice, we scale each kernel by a signal variance hyperparameter σ_f^2 , which governs the overall vertical variation in the function. This scaling is applied consistently across all kernel types and is discussed further in the 3.4.

From Table 1, we observe that the RBF, Rational Quadratic, Matern, Laplace, and Periodic kernels are all examples of stationary kernels (i.e depend on $|x - x'|$). Many of these—such as the RBF, Matern, Rational Quadratic and Laplace—exhibit “bell-shaped” structures: inputs x and x' that are close together yield high covariance, which then decays as the distance $\|x - x'\|$ increases. The Periodic kernel, while also stationary, has a unique structure. Instead of decaying monotonically with distance, it assigns high covariance to inputs that are separated by integer multiples of a fixed period p . This leads to a repeating pattern of similarity, making the kernel ideal for modeling functions which are periodic.

We have examined multiple kernel types and their properties. We will now briefly visualise and examine the effect of the lengthscale hyperparameter and the signal variance hyperparameter on the GP prior. We will use the RBF kernel as an example, but the same principles apply to other kernels.

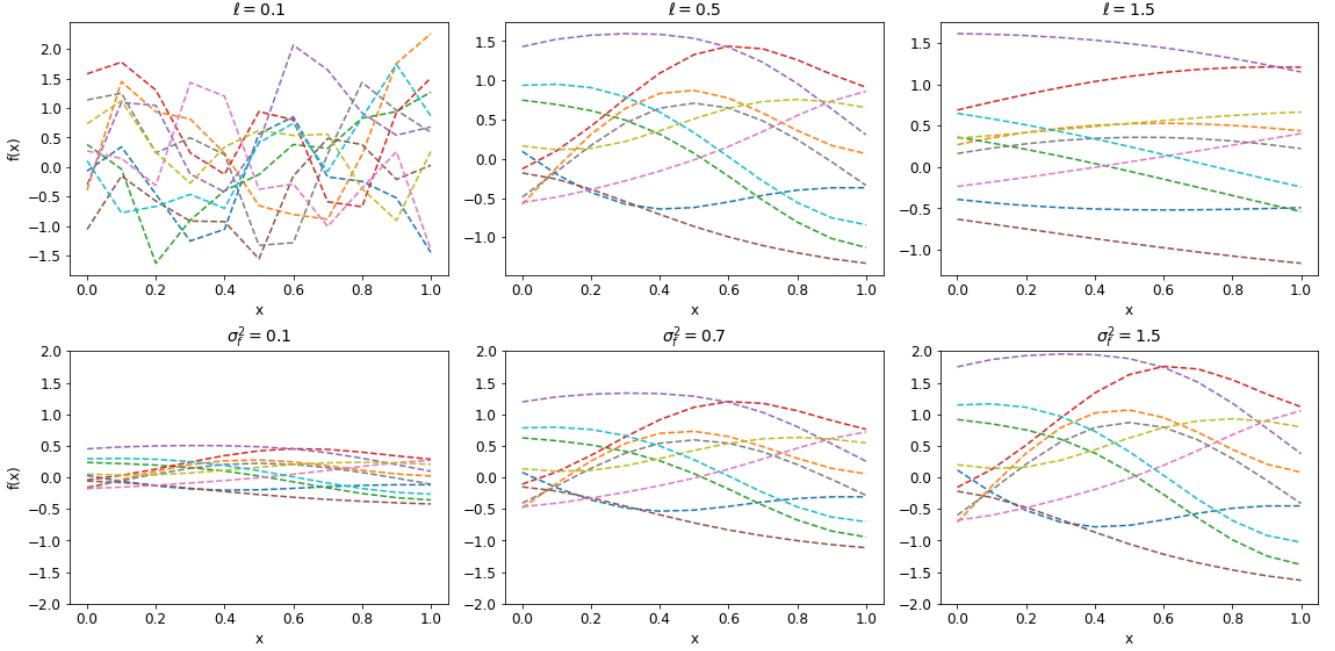


Figure 2: Sampling from the GP prior with mean 0 and covariance given by the RBF kernel. The first plot shows the effect of the lengthscale hyperparameter ℓ on the GP prior. We fix the signal variance to 1. The second plot shows the effect of the signal variance hyperparameter σ_f^2 on the GP prior. We fix the lengthscale to 0.5.

Sean: Comment on how each effects prior. Below is waffle from a while ago could potentially add

It is evident that the choice of kernel hyperparameters has a significant impact on the shape of the posterior mean and its uncertainty. The RBF kernel with noise is given by:

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right) + \sigma_n^2 I,$$

where ℓ controls the length scale, σ_f^2 represents the signal variance, and σ_n^2 denotes the noise level.

From Figure 2, we observe that a very small length scale $\ell \approx 0$ leads to overfitting: the GP closely tracks the training data, with sharp fluctuations between points. This occurs because the kernel's covariance rapidly decays with distance, making function values at different inputs nearly uncorrelated unless they are extremely close. As a result, the model fits the noise and exhibits a jagged appearance. In contrast, a large length scale such as $\ell \approx 3$ causes over-smoothing. Distant points remain highly correlated, producing a nearly linear mean function that fails to capture local variations. A moderate length scale ($\ell = 0.109$ in this case) provides a balance—producing a smooth yet flexible mean function that captures trends in the data without overfitting.

Regarding the signal variance σ_f^2 , we find that a very small value (e.g., $\sigma_f^2 \approx 0.01$) causes the GP mean to flatten, as the prior assumes the function to vary very little. As the variance increases (e.g., $\sigma_f^2 = 1.2$ and $\sigma_f^2 = 10$), the credible intervals widen slightly, but the overall mean function shape remains relatively stable. This suggests that posterior inference over σ_f^2 may exhibit considerable uncertainty, a point we explore further using MCMC sampling in Section 4.1.

Lastly, varying the noise level σ_n^2 primarily affects the width of the credible interval, while the mean function remains largely unchanged. As σ_n^2 increases, the model accounts for higher observation noise, and uncertainty in the predictions increases accordingly. This relationship appears approximately linear in the visualizations.

From previous section may be worth including Credible Interval Behavior: In Gaussian Process regression, the predictive variance at a new input x^* is given by

$$k(x^*, x^*) - k(x^*, X)(K_{XX}^{-1})k(X, x^*).$$

With an RBF kernel, $k(x^*, X)$ will be near-zero if x^* is farther than a few ℓ 's away from all training inputs X . For small ℓ , this situation happens frequently – any point outside a tiny neighborhood of training data effectively

has $k(x^*, X) \approx 0$, so its predictive variance is roughly

$$k(x^*, x^*) = \sigma_f^2$$

(the prior variance). Thus, the 95% credible interval reverts to roughly $\pm 2\sigma_f$ in large gaps or outside the data range. For large ℓ , on the other hand, $k(x^*, X)$ remains sizable over a much broader range, reducing the variance term. The uncertainty only approaches the prior level far outside the data (several ℓ 's away). In effect, a large ℓ flattens the covariance function so that the GP “remembers” the training set far out, maintaining narrower error bars over a wider domain.

Adding Data: Prior to Posterior

We have discussed how our prior distribution is dependent on the choice of kernel and the hyper-parameters of said kernel. In this section we will discuss how we can update our prior distribution 18 to achieve our new prediction distribution distribution from which we can make inferences. One of the key strengths of Gaussian Processes is that, given observations at training inputs X and setting kernel hyper-parameters θ we can make predictive inferences about the function value at any new test location x_* . By applying the standard conditional Gaussian formulas (see appendix [Must clean up this derivation in appendix A](#) for the full derivation), the posterior distribution of $f(x_*)$ given $\{X, f(X)\}$ is Gaussian and given by:

$$p(f(x_*) | f(X), X, X_*, \theta) = \mathcal{N}(m(x_*), \sigma^2(x_*)), \quad (20)$$

where

$$m(x_*) = \mu(x_*) + k(x_*, X) k(X, X)^{-1} [f(X) - \mu(X)], \quad (21)$$

$$\sigma^2(x_*) = k(x_*, x_*) - k(x_*, X) k(X, X)^{-1} k(X, x_*). \quad (22)$$

In these expressions:

- $\mu(\cdot)$ is the mean function (We take this to be zero since we centre the data prior to prediction),
- $k(X, X)$ is the covariance matrix among the observed training points,
- $k(x_*, X)$ is the vector of cross-covariances between the test point x_* and the training inputs,
- $k(x_*, x_*)$ is the prior variance at the test point itself.

We now have an analytic function that can be evaluated to find the mean function value and variance at any input point. This is a very useful property that Gaussian Processes possess. Figure 4 shows how we update our distributions based on the training points. We initially plot samples taken from the prior distribution (Equation 18). After conditioning this distribution on one training point and obtaining the new predictive posterior distribution (Equation 20), we see that the predictive mean passes exactly through the training point, has small variance around it, and then fans out farther away. Conditioning on two training points at opposite ends of our input domain creates an ellipse-shaped credible interval whose largest radius appears midway between the two training points. With more training points, the predictions align progressively closer to the true function and the variance becomes smaller.

Note that this example uses a one-dimensional, noise-free function and a basic RBF kernel with fixed hyperparameters ($\ell = 1$, $\sigma^2 = 0.5$). In practice the choice of Kernel function plays a pivotal role in the assumptions we make about the shape and general behaviour of our model. In the next section we examine the different kernel choices available and the assumptions that each kernel encodes about our function structure, such as smoothness and periodicity.

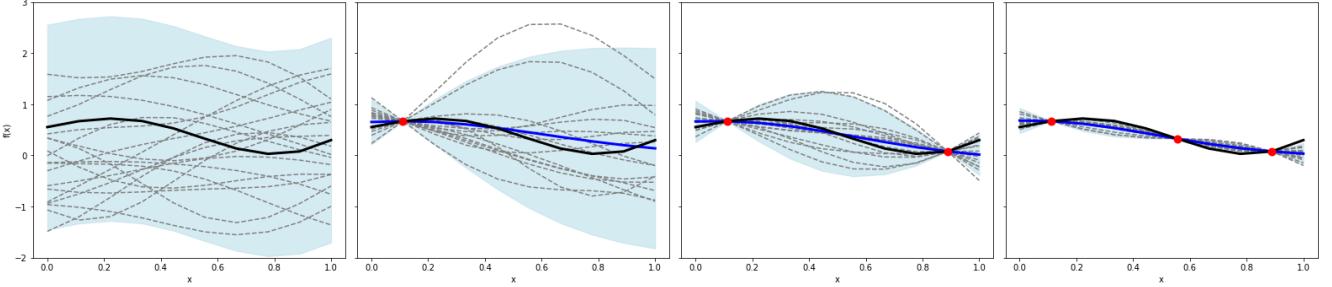


Figure 3: Prior Distribution

Figure 4: 1D Gaussian Process Regression: Prior to Posterior. This sequence shows how the GP prior transforms into a posterior as more data points are added. The RBF kernel was used with hyper-parameters: $l = 1$, $\sigma^2 = 0.5$. The black line represents the true function. The blue is the mean of each posterior distribution. The light blue shaded region is the credible interval and the grey lines are the samples drawn from each posterior/prior.

Sean: can add a potential link to animation here illustrating the nice distribution formed at each point

3.4 Handling Noise in our Data

So far, our discussion has assumed noise-free observations. However, real-world data is rarely clean, measurements often include some form of uncertainty. To make our Gaussian Process models more applicable to this real-world data, we now explore how to incorporate noise into the GP framework.

We assume that each observation includes the true function value plus Gaussian noise:

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_{n,i}^2), \quad (23)$$

where $\sigma_{n,i}^2$ is the noise variance associated with input x_i . This allows us to model both homoscedastic and heteroscedastic noise under a unified notation.

Under this model, we have the following Gaussian assumptions:

$$f \sim \mathcal{N}(0, K), \quad (\text{prior over the true function}) \quad (24)$$

$$y \sim \mathcal{N}(0, K + \Sigma), \quad (\text{distribution over noisy observations}), \quad (25)$$

where $\Sigma = \text{diag}(\sigma_{n,1}^2, \sigma_{n,2}^2, \dots, \sigma_{n,n}^2)$ is the noise covariance matrix.

This now updates our previous posterior mean and variance (eq: 21 and 22) to a revised posterior:

$$P(f_*|X, X_*, \theta, y) = \sim \mathcal{N}(m(f_*), \text{Var}(f_*)) \quad (26)$$

$$m(f_*) = K_*^T (K_y)^{-1} y, \quad (27)$$

$$\text{Var}(f_*) = K_{**} - K_*^T (K_y)^{-1} K_*. \quad (28)$$

where K_y depends on how we handle our noise. There are three main cases of how we handle our noise

Homoscedastic Noise

In the homoscedastic case, we assume that all observations have the same noise level, meaning the noise variance is constant across the dataset:

$$\sigma_{n,i}^2 = \sigma_n^2 \quad \forall i.$$

This simplifies the noise covariance matrix Σ to a scalar multiple of the identity matrix:

$$\Sigma = \sigma_n^2 I.$$

The total covariance matrix of the observed data becomes:

$$K(X, X) + \sigma_n^2 I = \begin{bmatrix} k(x_1, x_1) + \sigma_n^2 & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) + \sigma_n^2 \end{bmatrix}.$$

Our prior distribution now becomes:

$$y \sim \mathcal{N}(0, K + \sigma_n^2 I) \quad (29)$$

where σ_n^2 is a new parameter which effects the shape of the distribution. In figure 2 we explored how the internal kernel hyper-parameters effect the shape of the samples from our prior distribution. We now examin how the noise (i.e σ_n^2) effects samples from our prior distribution from 26. Our predictive distribution remains as in 26

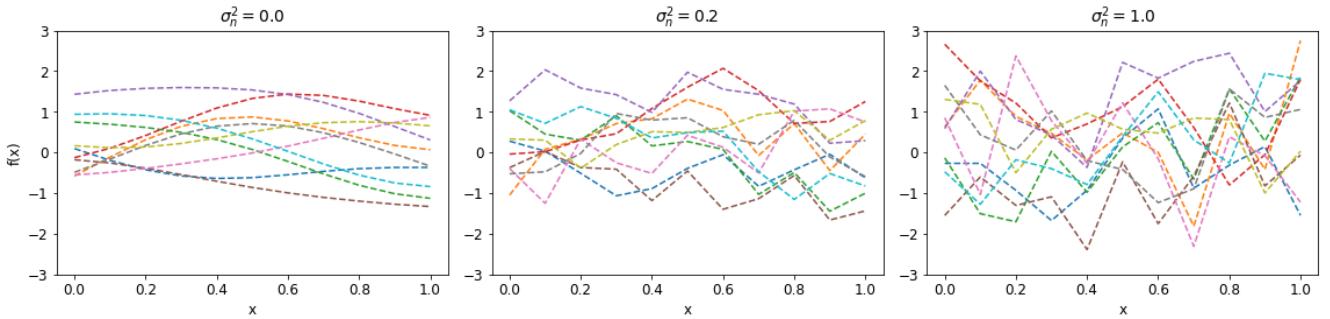


Figure 5: Sampling from the GP prior with mean 0 and covariance given by the RBF kernel. The plot shows the effect of the noise hyperparameter σ_n^2 on the GP prior. We fix the signal variance to 1 and length scale to 0.5.

Our posterior distribution is as in eqn 26 with $K_y = K(X, X) + \sigma_n^2 I$ for some constant σ_n

Heteroscedastic Noise

Known Noise In this case, the noise variance changes across the input space—some observations are noisier than others. If we know the individual noise variances σ_i^2 for each training input x_i , we incorporate them by adding a diagonal noise matrix to the kernel:

$$K(X, X) + \Sigma = \begin{bmatrix} k(x_1, x_1) + \sigma_1^2 & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) + \sigma_n^2 \end{bmatrix},$$

where $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$. In this case the noise is not found as a hyper-parameter but instead just added to the diagonal of the covariance matrix. Our posterior distribution is as in eqn 26 with $K_y = K(X, X) + \Sigma$ where Σ is the known noise of our data.

Learning Noise over the Input Space If the noise variance is unknown but varies across the input space, we can model it as a function. This is done by building a kernel that captures both smooth, global trends and rough, local fluctuations. In practice, this means building an additive kernel made up of sub-kernels. For example, as seen in Table 1, some kernels like the Matern, Laplacian, or Rational Quadratic capture local variations well (interpreted as noise), while others like the RBF capture broader, smoother structure. By combining these, we can allow one kernel component to model the general structure of the function, and the other to model the heteroscedastic noise behavior. Our posterior distribution is as in eqn 26 with :

$$K_y = \theta_1 K_1(X, X) + \theta_2 K_2(X, X), \quad (30)$$

where K_1 and K_2 are distinct kernels chosen to capture different aspects of the data. The coefficients θ_1 and θ_2 are parameters that control the relative contribution of each kernel component.

Monte Carlo Sampling of Noise

This technique can be applied to both homoscedastic and heteroscedastic noise settings. Instead of modifying the kernel matrix analytically, we simulate the effect of noise by generating multiple noisy realizations of the observed outputs. We assume the observation noise is Gaussian:

$$\epsilon_i \sim \mathcal{N}(0, \sigma_i^2),$$

and therefore, the observed data y is a noisy version of the true latent function values:

$$y = f + \epsilon, \quad \text{with } y \sim \mathcal{N}(f, \Sigma),$$

where $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. To simulate this, we generate M noisy samples of the observations:

$$y^{(s)} = y + \epsilon^{(s)}, \quad \epsilon^{(s)} \sim \mathcal{N}(0, \Sigma) \quad (31)$$

For each sampled dataset $y^{(s)}$, we compute an independent GP posterior:

$$p(f_* | X, X_*, \theta, y^{(s)}). \quad (32)$$

To obtain the final predictive distribution, we marginalize over these sampled posteriors:

$$p(f_* | X, X_*, \theta, y) = \int p(f_* | X, X_*, \theta, y^{(s)}) p(y^{(s)} | y) dy^{(s)}. \quad (33)$$

This integral is intractable in closed form, so we approximate it using Monte Carlo integration:

$$p(f_* | X, X_*, \theta, y) \approx \frac{1}{M} \sum_{s=1}^M p(f_* | X, X_*, \theta, y^{(s)}). \quad (34)$$

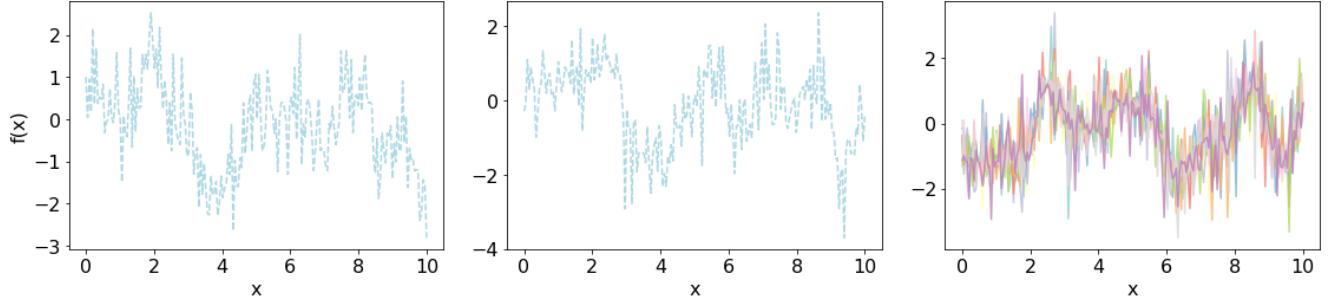


Figure 6: Samples drawn from a zero-mean Gaussian Process prior with varying noise assumptions. **Left:** Homoscedastic noise, where a constant noise variance $\sigma_n^2 = 0.5$ is added uniformly across all inputs. **Middle:** Heteroscedastic noise, where individual noise variances σ_i^2 are known and drawn from $\mathcal{N}(0, 0.5)$, resulting in a diagonal noise covariance. **Right:** Monte Carlo sampling of noisy observations, where multiple noisy realizations are generated from $\mathcal{N}(f(x), \epsilon^2)$

Sean: Own words

We can observe in Figure 6 how different noise assumptions influence samples drawn from the Gaussian Process prior. In the homoscedastic case (left), the output exhibits consistent fluctuations across the entire domain due to a constant noise variance applied uniformly to all inputs. In contrast, the heteroscedastic case (middle) introduces input-dependent noise, resulting in regions of smoothness interspersed with abrupt variations—reflecting the fact that each output has its own associated noise level. Finally, in the Monte Carlo noise sampling approach (right), we generate multiple realizations by perturbing the data with different noise samples. This highlights the diversity of plausible functions consistent with the observed data and captures the full range of uncertainty introduced by noisy observations.

3.5 Hyper-parameters

Until now, all predictive distributions such as Equations 26 and 20 have been conditioned on fixed kernel hyperparameters. As demonstrated in Figures 2 and 5, these hyperparameters have a significant influence on the structure and behaviour of the Gaussian Process, shaping both the prior and posterior distributions.

In practice, these hyperparameters are not known a priori and must be inferred from the data. To do so, we seek the set of hyperparameters that best explain the observed data by maximising the *log marginal likelihood*, a method detailed in Chapter 5 of [6].

From Equation 26, we have:

$$y \sim \mathcal{N}(0, K + \sigma_n^2 I)$$

where K is the kernel matrix computed from the training inputs X , and σ_n^2 is the noise variance.

This implies that the marginal likelihood (i.e., the probability of the observed outputs y given the inputs X and hyperparameters θ) is given by the multivariate Gaussian density:

$$p(y | X, \theta) = \frac{1}{(2\pi)^{n/2} |K_y|^{1/2}} \exp\left(-\frac{1}{2} y^\top K_y^{-1} y\right)$$

where $K_y = K + \sigma_n^2 I$.

Taking the logarithm of this expression yields the *log marginal likelihood*:

$$\log p(y | X, \theta) = -\frac{1}{2} y^\top K_y^{-1} y - \frac{1}{2} \log |K_y| - \frac{n}{2} \log 2\pi \quad (35)$$

Our goal is to maximise this log marginal likelihood with respect to the hyperparameters θ , which typically includes the kernel lengthscale, signal variance, and noise variance. Once optimal values are found, we can use them to make accurate posterior predictions.

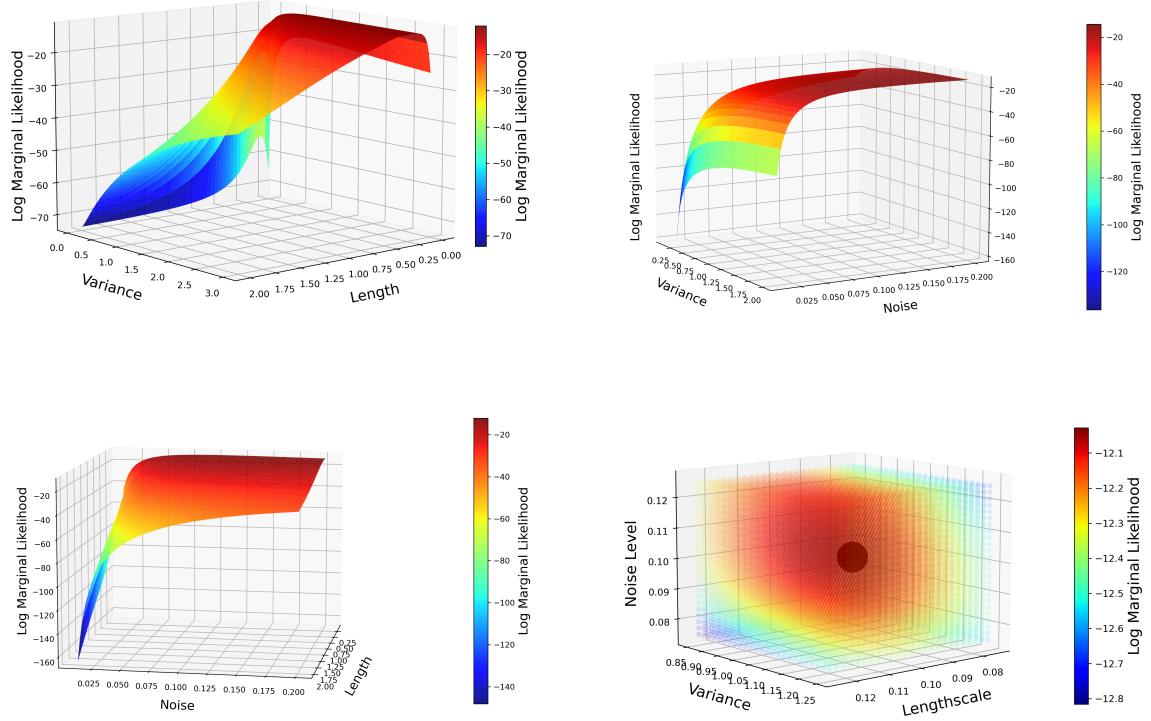


Figure 7: For a GPR with noise we are forced to optimise a length hyper-parameter, a variance hyper-parameter and a noise hyper-parameter. Here we have kept one parameter constant on each graph and compared the log likelihood space varying the other two parameters. In graph one the noise is set at 0.1, in graph 2 the length is set at 0.5, in graph 3 the variance is set at 1.5. In the final graph we plot a 3-dimensional scatter plot and illustrate the point estimate given by the optimisation algorithm as the black dot. This point is located at $(\sigma^2 = 1.16, l = 0.109, \text{noise} = 0.105)$

3.5.1 Cross Validation

While maximising our parameters with respect to the log likelihood is common practice [ch 5, [6]] details how we may overfit our model to the training data. To account for this we use cross validation where we split our data into disjoint sets. One set is used to train our model (find the hyper-parameters) and the other set is used to monitor performance. If we split our data into K sets every iteration we train on $K-1$ sets, and validate on the k th set. We iterate this until all sets have been used to validate our model.

Sean: Have to think about this model selection is what I am talking about here. Maybe leave out cross val here because I am not doing cross val for hyper-parameters

4 Quantifying Uncertainty and Evaluating Model Accuracy

4.1 Monte Carlo Markov Chain (MCMC): A Posterior over Hyperparameters

Why MCMC and the Mathematics Behind It

We have discussed the role of kernels and how their hyperparameters affect the posterior distribution. From Figures 2 and 7, we observed that when optimising hyperparameters, there is not necessarily a single "best" solution, but rather a region of valid values that explain the data well.

Up to this point, the models we considered have relied on point estimates—selecting the hyperparameters that maximise the log marginal likelihood and using them for predictions. However, this approach ignores the uncertainty between hyperparameters that yield similarly high log-likelihood scores. To address this uncertainty, we adopt a Bayesian perspective by using Markov Chain Monte Carlo (MCMC) methods to sample from the posterior distribution over hyperparameters. This approach allows us to move beyond point estimates and instead capture a full distribution that reflects uncertainty and variability in the hyperparameters. By leveraging this posterior distribution, we gain a more comprehensive understanding of the model’s behaviour. It allows for more robust predictions, better uncertainty quantification, and improved decision-making—especially in cases where multiple hyperparameter settings are plausible.

When getting a point-estimate we maximised the log likelihood from equation 35. Now we want to build a distribution over the hyper-parameters. We have that:

$$p(\theta | \mathbf{y}, X) = \frac{p(\mathbf{y} | X, \theta) p(\theta)}{p(\mathbf{y} | X)}$$

where:

- $p(\mathbf{y} | X, \theta)$ is the likelihood,
- $p(\theta)$ is the prior over hyperparameters,
- $p(\mathbf{y} | X)$ is the marginal likelihood, acting as a normalising constant.

Since $p(\mathbf{y} | X)$ is often intractable, we sample from the unnormalised posterior using MCMC:

$$p(\theta | \mathbf{y}, X) \propto p(\mathbf{y} | X, \theta) p(\theta)$$

Using MCMC, we generate samples $\{\theta^{(s)}\}_{s=1}^S \sim p(\theta | \mathbf{y}, X)$ from this posterior.

Implementing MCMC

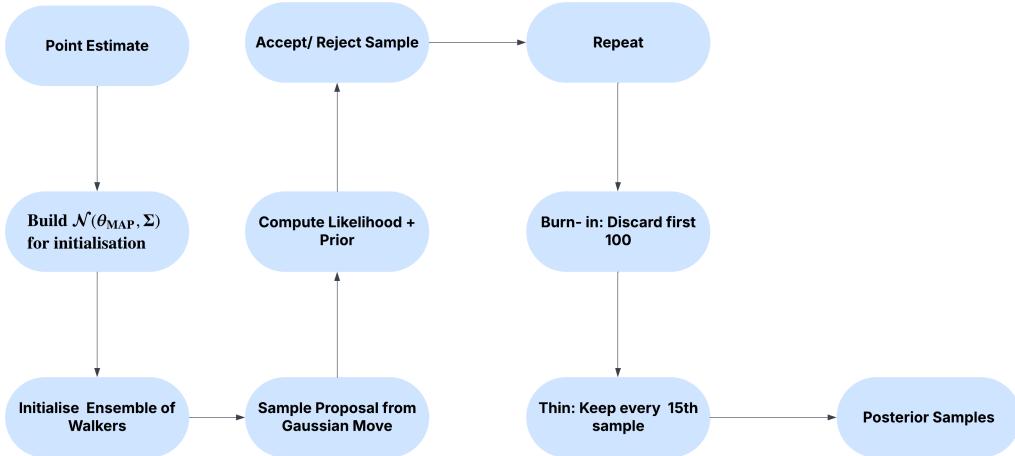


Figure 8: Overview of the MCMC sampling procedure for Gaussian Process hyperparameter inference. This pipeline samples from the posterior $p(\theta | \mathbf{y}, X)$ using a Metropolis-Hastings Gaussian proposal and an ensemble of walkers.

Before starting MCMC, we must decide which model we want to build the hyperparameter posterior for. This involves selecting one of the six kernels outlined in Section 3.3, along with one of the three noise-modelling approaches described in Section 3.4. Once this model structure is fixed, we obtain initial point estimates for the hyperparameters by maximising the log marginal likelihood, as discussed in Section 3.5.

We then construct a multivariate normal distribution centred at this point estimate and sample from it to initialise each walker. From there, the walkers explore the hyperparameter space using a Gaussian proposal distribution with a specified covariance. At each step, we compute the sum of the log likelihood and log prior. The proposed sample is then accepted or rejected using the Metropolis-Hastings criterion Could give more detail here. This process is repeated for a fixed number of steps to generate a large set of samples.

To ensure convergence and sample independence, we discard the first 100 samples from each walker as burn-in and apply a thinning factor of 15—retaining every 15th sample. The resulting collection of samples forms our posterior distribution over hyperparameters, which we visualise using a kernel density estimate (KDE).

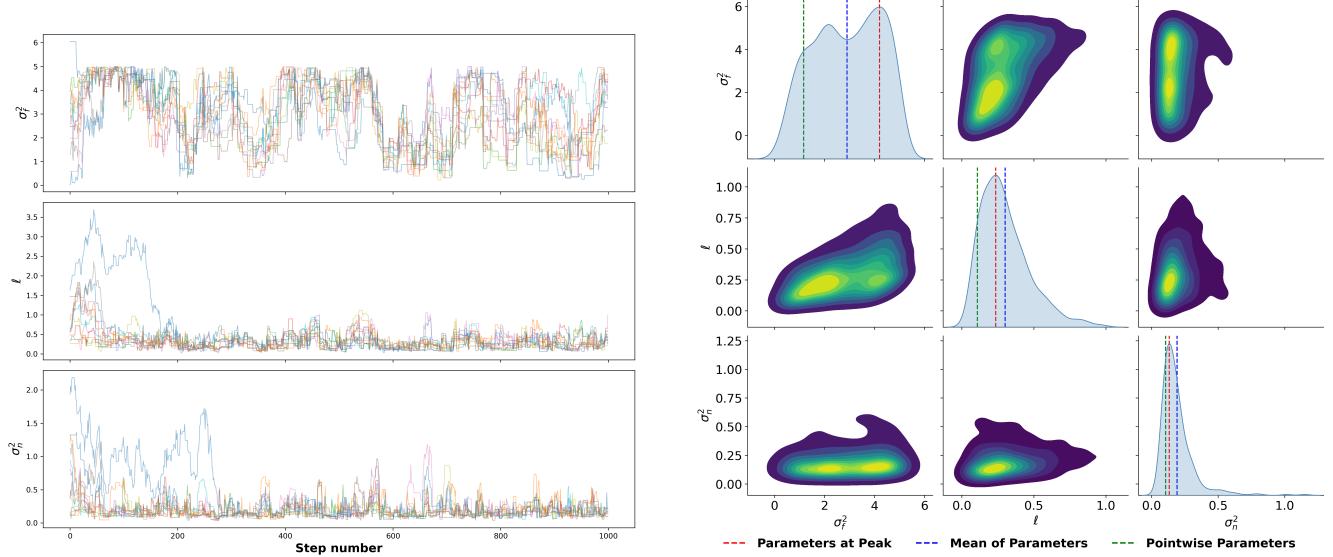


Figure 9: This MCMC is run using a gpr with a RBF Kernel with White Kernel for noise. Left: A plot of the space sampled by each walker in my MCMC run (burnin=100 and thin=15). Right: A plot of the distribution of the hyper-parameters built using a Gaussian KDE from the MCMC samples.

We can see from the distributions in 9 that the variance hyper-parameter is almost bi-modal (i.e has two significant peaks). For our bi-modal hyper-parameter we can see that the mean parameter, peak parameter and point estimate parameter differ significantly. In this scenario any singular point estimate will lose a lot of information about the distribution particularly for the variance hyper-parameter. To resolve this we can instead of picking a singular point build the hyper-parameter uncertainty into our final predictive distribution.

We previously had our predictive distribution as :

$$p(f_* | \mathbf{y}, X, X_*, \theta),$$

where predictions were made conditional on a fixed set of hyperparameters θ . Now, we marginalise over the posterior distribution of θ to account for hyperparameter uncertainty:

$$p(f_* | \mathbf{y}, X, X_*) = \int p(f_* | \mathbf{y}, X, X_*, \theta) p(\theta | \mathbf{y}, X) d\theta.$$

Since this integral is analytically intractable, we approximate it using the MCMC samples $\{\theta^{(s)}\}_{s=1}^S$:

$$p(f_* | \mathbf{y}, X, X_*) \approx \frac{1}{S} \sum_{s=1}^S p(f_* | \mathbf{y}, X, X_*, \theta^{(s)}).$$

This entails constructing a full posterior predictive distribution for each set of sampled hyperparameters and then averaging across all predictions. In practice, we compute the final predictive mean and variance using the law of total variance:

$$\mathbf{E}[f_*] \approx \frac{1}{S} \sum_{s=1}^S \mu^{(s)}(f_*),$$

$$\text{Var}[f_*] \approx \frac{1}{S} \sum_{s=1}^S \left[\sigma^{2(s)}(f_*) + \left(\mu^{(s)}(f_*) \right)^2 \right] - (\mathbf{E}[f_*])^2,$$

where $\mu^{(s)}(f_*)$ and $\sigma^{2(s)}(f_*)$ are the predictive mean and variance obtained from the s -th sampled hyperparameter configuration $\theta^{(s)}$.

This procedure yields a predictive distribution that fully reflects both the uncertainty in the data and the uncertainty in the model's hyperparameters

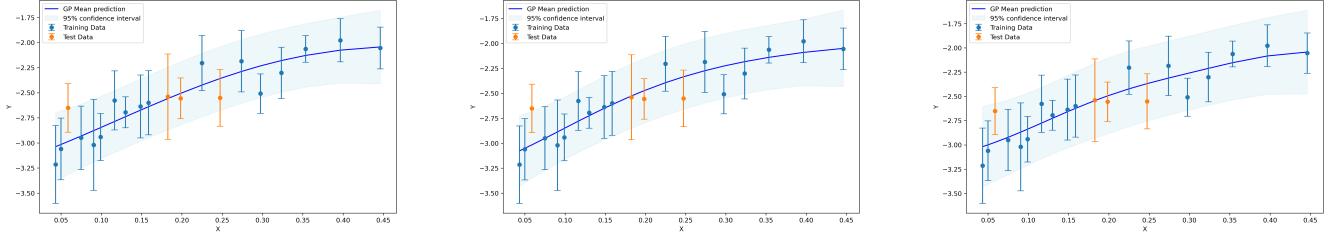


Figure 10: Left: This is the GPR plotted with the mean hyper-parameters Middle: This is the GPR plotted with the peak hyper-parameters Right: This is the full predictive distribution over all the hyper-parameter samples

4.2 Model Evaluation Metrics

In figure 2, we made a graphical representation of four out of six metrics used to evaluate our model's accuracy. [6] discusses how these metrics provide a balanced assessment of model performance. The Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) measure the average deviation of predictions from the true values, with RMSE penalizing larger errors more heavily. The coefficient of determination R^2 quantifies how well the model predicts relative to the mean of the test set. It is computed as 1 minus the ratio of the squared residuals to the total variance. A value closer to 1 indicates better predictive performance. The adjusted R^2 (\bar{R}^2) accounts for model complexity by penalizing excessive predictor variables, preventing overfitting. The Figure of Merit (FOM) evaluates the ratio of a point's prediction error to its associated standard deviation. A FOM near 1 is ideal, indicating that the model's uncertainty estimates are well-calibrated. A FOM $\ll 1$ suggests an overly conservative model with large uncertainty, while a FOM $\gg 1$ may indicate overconfidence, failing to capture true variability. The Pearson correlation coefficient measures the linear relationship between predictions and true values. A correlation of 1 (-1) signifies a perfect positive (negative) linear relationship, whereas a correlation of 0 indicates no linear association.

Metric Name	RMSE	R^2	FOM	Pearson Coefficient
Formula	$\sqrt{\frac{1}{N} \sum (y_i - \hat{y}_i)^2}$	$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$	$\frac{\text{RMSE}}{\sigma}$	$\frac{\text{cov}(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}}$
Visual Illustration				

Table 2: Comparison of different performance metrics used in evaluating models. RMSE, R^2 , FOM, and the Pearson Coefficient are included. MAE is similar to RMSE but without squaring errors. Adjusted R^2 accounts for the number of predictors and is slightly modified from R^2 . The actual metrics for each graph are: RMSE = 0.2, R^2 = 0.6, FOM = 1.09, Pearson correlation = 0.8.

5 Multi-Dimensional Gaussian Process Regression

- Here we want to bring up how we can have a different length parameter for each dimension of the data.
- Get a few plots to demonstrate this maybe.
- Important to note noise modelling for all dimensions is the same since just adding to the diagonal of the covariance matrix.
- Mention how hyper-parameter optimisatino changes in the larger parameter space

6 Method

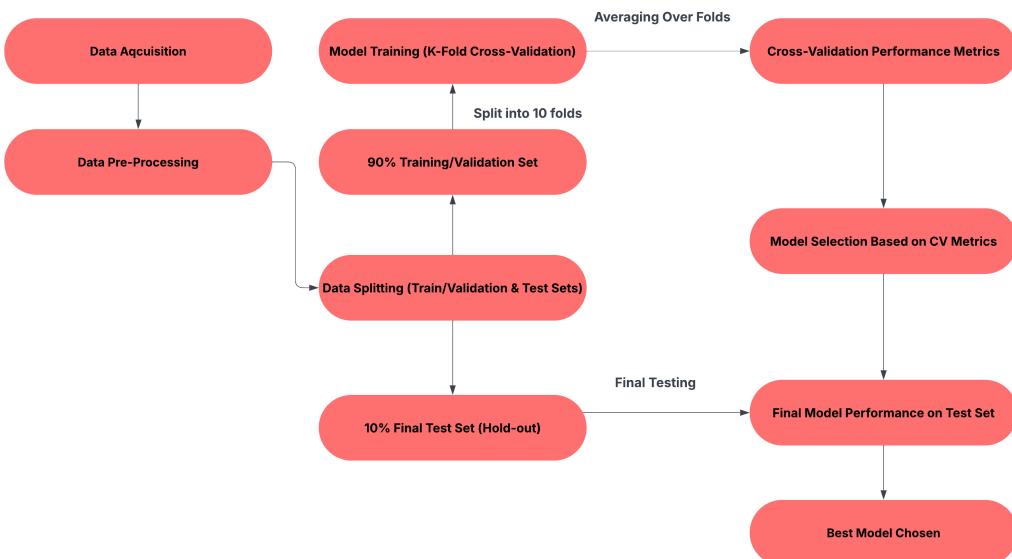


Figure 11: My Process Flowchart

6.1 Data Description

- Dimensionality choices (4D vs 7D) - Summary of input features - How the data was generated or collected (e.g., from simulations or experiments)

Kernel Choices

I built my selection of models using 5 kernels discussed in section 3.3 these kernels are the Radial Basis Function (RBF), the Rational Quadratic Function (RQF), the Matern Kernel (Mat), the periodic kernel (Exp sine) and the laplace kernel (lapl).

Noise approach

I used each noise type discussed in section 3.4, this includes treating noise as a hyper-parameter (homoscedastic noise), using our known noise (heteroscedastic), using combined kernels to model the noise (heteroscedastic), using monte carlo sampling of the noise (heteroscedastic).

Hyper parameters

For all my models I maximised the marginal log-likelihood for parameter selection as discussed in 3.5. Then after finalising models I also ran Monte Carlo Markov Chain simulations to build a posterior over the hyper-parameters and understand the uncertainty/ distribution associated with each hyper-parameter.

Metrics Calculation

In order to compare these models effectively I evaluated each model using the metrics discussed in section 4.2 and then compared the metrics.

Cross Validation

To ensure model's generalised well and weren't biased to a certain training set I used 10 fold cross validation and on each fold calculated the metrics for comparison.

7 Results

7.1 Cross-Validation Performance

- Plots of metrics per model type - Ranking tables or summaries

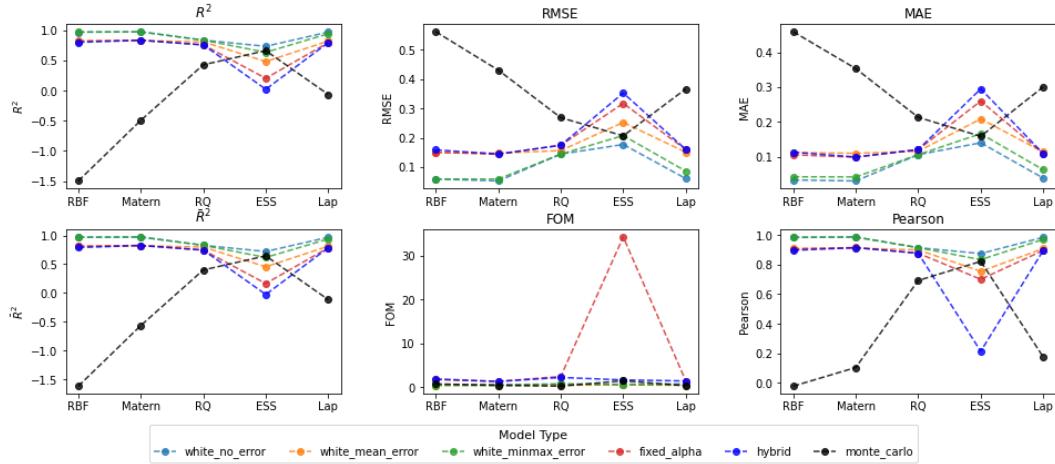
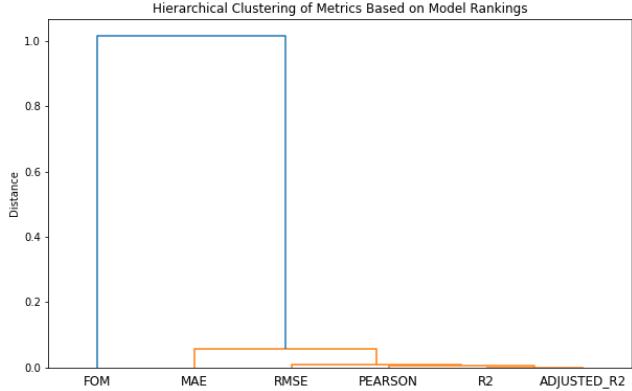
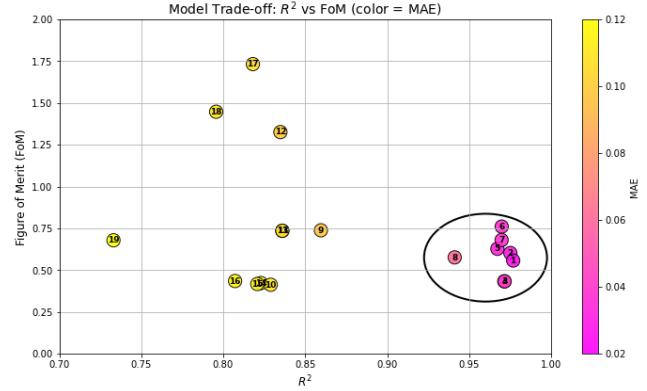


Figure 12: Comparing how similar results the metrics give in a dendrogram



(a) Results from cross-validation.



(b) Comparing clustered model performance by R^2 and FoM.

Figure 13: Cross-validation metric trends and a performance trade-off visualization.

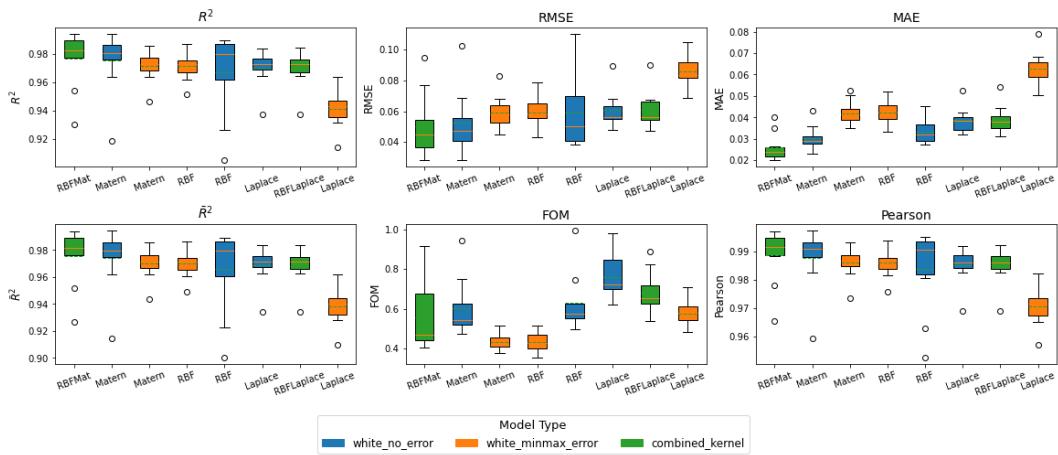


Figure 14: Seeing how the best models performed over different clusters

7.2 Training on 90% of Data

- Final model evaluations on held-out 10- Error bars, confidence intervals, diagnostics

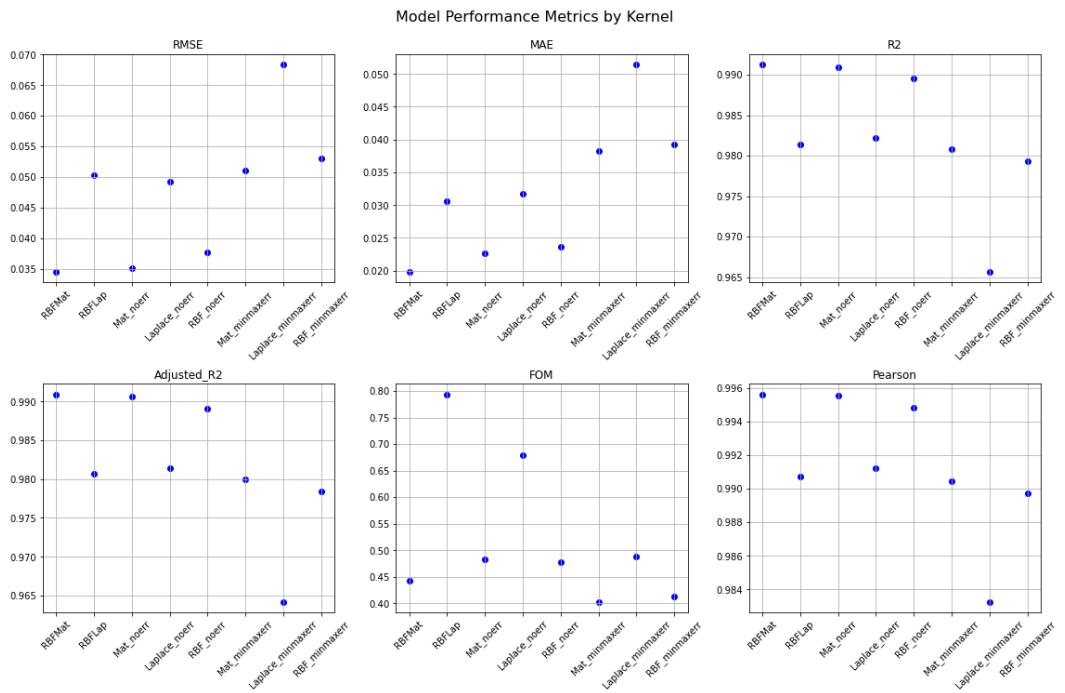


Figure 15: Metrics of best 8

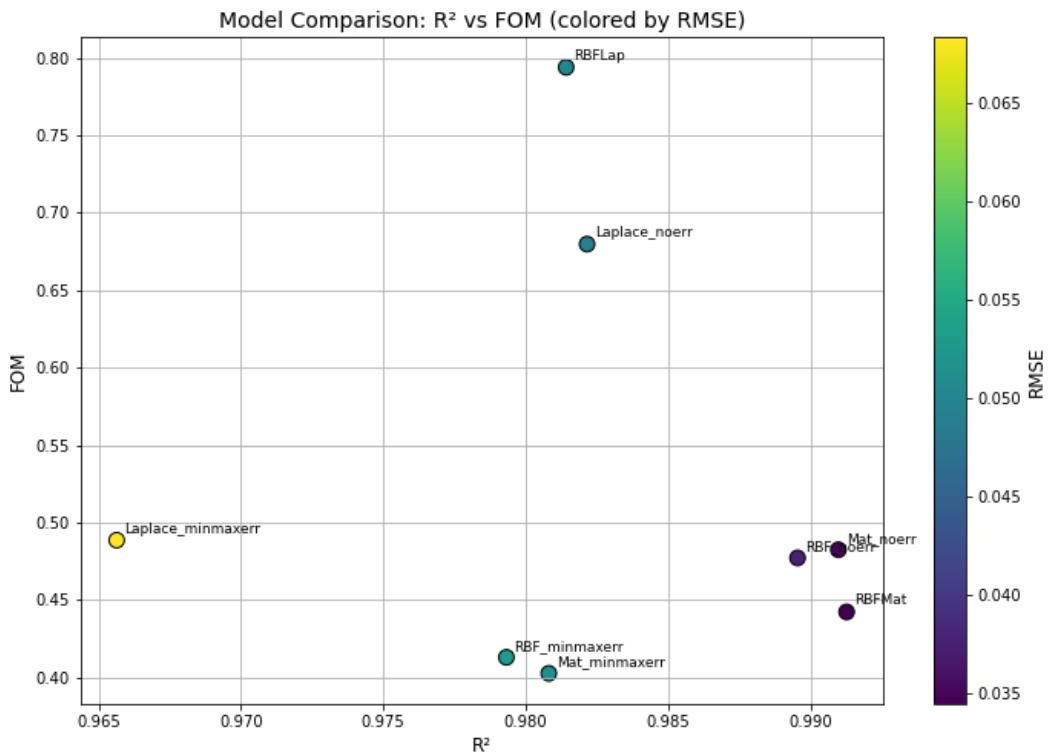


Figure 16: Comparing metrics

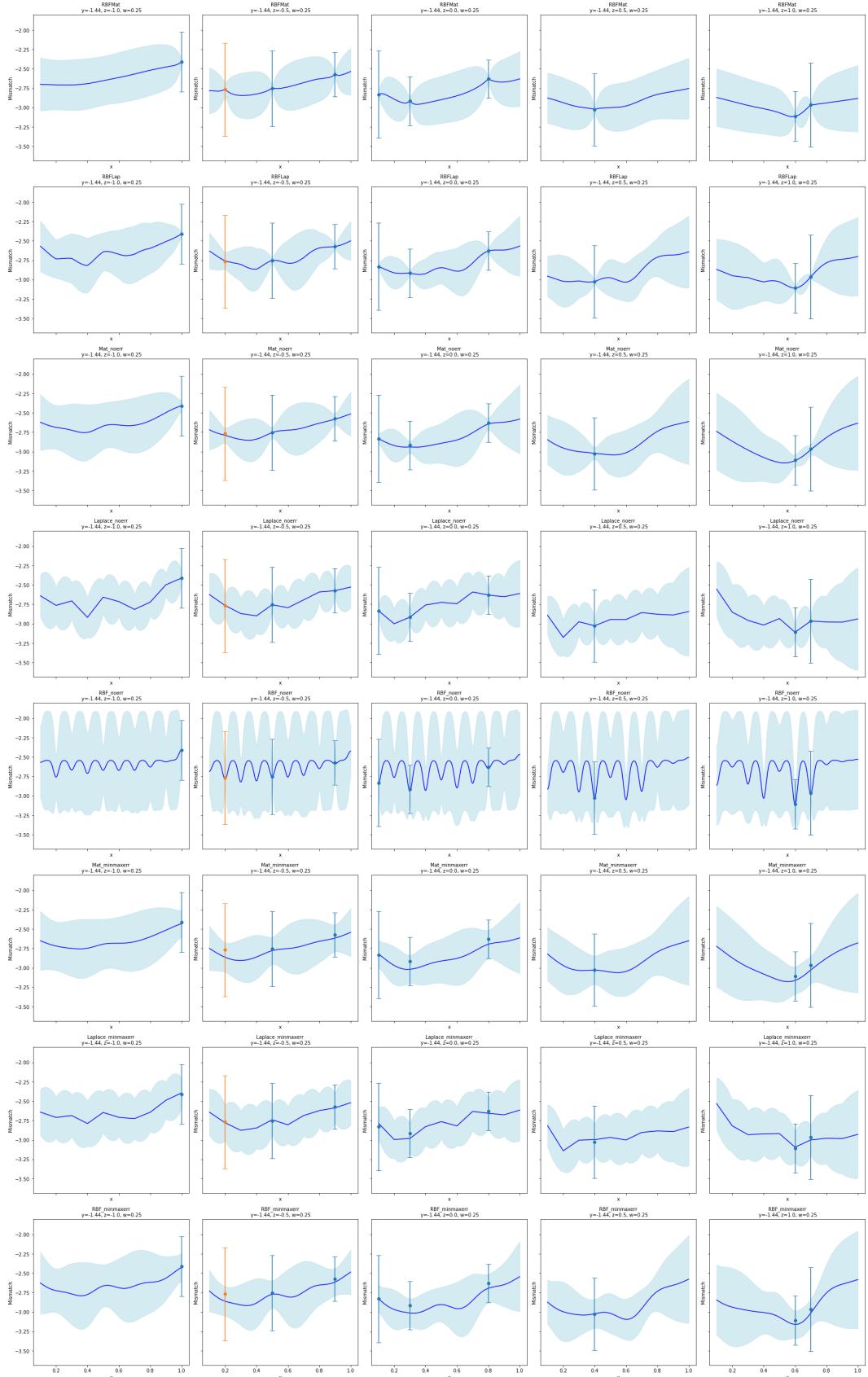


Figure 17: All 8 gps with cutting their y-axis

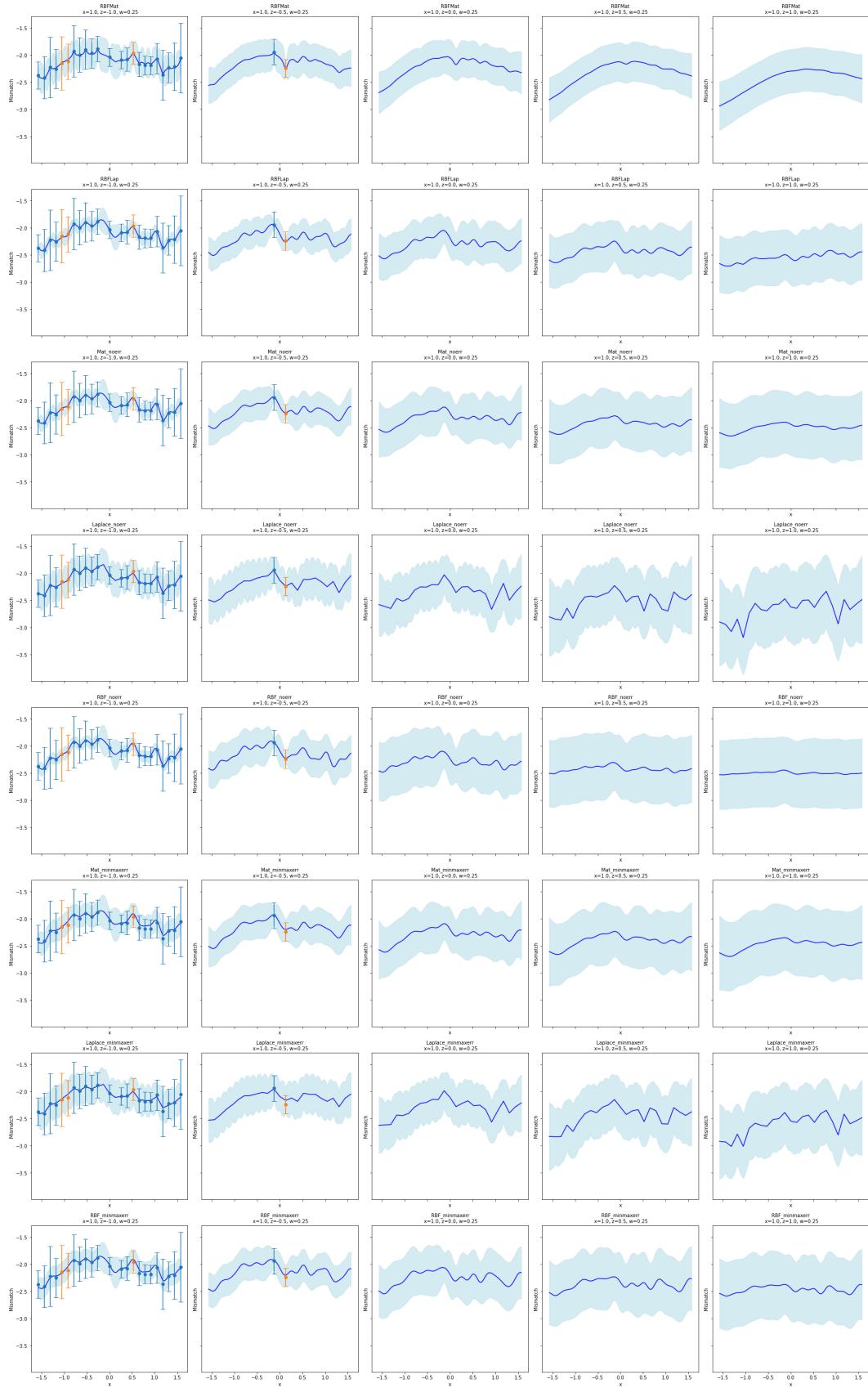


Figure 18: All 8 gps with cutting their x-axis

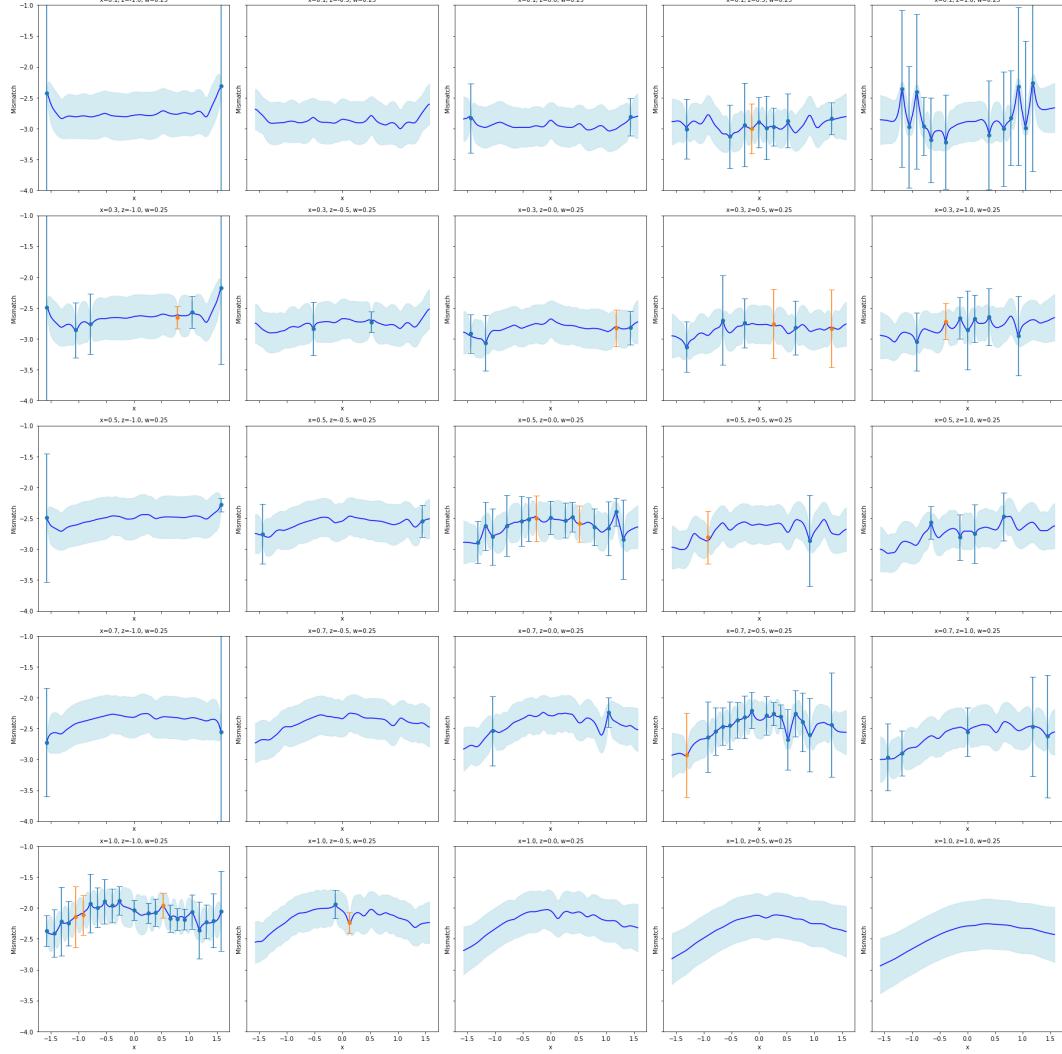


Figure 19: Examining my chosen model RBF Matern kernel

7.3 Top Model Selection and Analysis

- Your top 3 models - Summary of their performance - MCMC trace plots or hyperparameter uncertainty discussion

A Appendix A: Derivation of Predictive Distribution

Using Bayes' Theorem applied to continuous probabilities, we have:

$$p(f_*|f) = \frac{p(f_*, f)}{p(f)}.$$

we have:

$$p(f_*, f) = \frac{1}{2\pi\sqrt{|\mathbf{C}|}} \exp\left(-\frac{1}{2} \begin{bmatrix} f \\ f_* \end{bmatrix}^T \mathbf{C}^{-1} \begin{bmatrix} f \\ f_* \end{bmatrix}\right)$$

and

$$p(f) = \frac{1}{\sqrt{2\pi}|K|^{1/2}} \exp\left(-\frac{1}{2} f^T K^{-1} f\right).$$

Legend

- **Covariance matrices:**

- $K = K(X, X)$: Covariance matrix of the training inputs.
- $K_{**} = K(X_*, X_*)$: Covariance matrix of the test inputs.
- $K_* = K(X, X_*) = K(X_*, X)^\top$: Cross-covariance between training and test inputs.

- **Joint covariance matrix:**

$$C = \begin{bmatrix} K & K_* \\ K_*^\top & K_{**} \end{bmatrix}$$

- **Determinant of C:**

$$|C| = KK_{**} - K_*K_*^\top$$

- **Inverse of C:**

$$C^{-1} = \frac{1}{|C|} \begin{bmatrix} K_{**} & -K_* \\ -K_*^\top & K \end{bmatrix}$$

- **Mean functions:**

$$m(X) = m(X_*) = 0$$

B Appendix B: Kernel Formulas

Radial Basis Function (RBF) Kernel

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

This kernel assumes smooth and infinitely differentiable functions, modeling local variations.

Rational Quadratic Kernel

$$k(x, x') = \sigma_f^2 \left(1 + \frac{(x - x')^2}{2\alpha\ell^2}\right)^{-\alpha}$$

This kernel can be seen as a scale mixture of RBF kernels, allowing for multi-scale behavior.

Periodic Kernel

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\frac{\pi(x - x')}{p}\right)\right)$$

This kernel models repeating structures with period p .

Matern Kernel

$$k(x, x') = \sigma_f^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|x - x'|}{\ell}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}|x - x'|}{\ell}\right)$$

The Matern kernel allows for controlling the smoothness of functions via the parameter ν .

Laplace (Exponential) Kernel

$$k(x, x') = \sigma_f^2 \exp(-\gamma|x - x'|)$$

Equivalent to the Matern kernel with $\nu = \frac{1}{2}$, this kernel models rougher functions.

Linear (Dot-Product) Kernel

$$k(x, x') = \sigma_b^2 + x^\top x'$$

This kernel grows with the similarity (inner product) between inputs, and it allows the function to vary globally. Since it depends directly on the values of x and x' , not just their difference, it is non-stationary. It is particularly useful for modeling linear trends.

We obtain the posterior predictive distribution:

$$p(f_*|f) = \frac{1}{(2\pi)^{n/2}|K_{**} - K_*^T K^{-1} K_*|^{1/2}} \exp\left(-\frac{1}{2}(f_* - K_*^T K^{-1} f)^T (K_{**} - K_*^T K^{-1} K_*)^{-1} (f_* - K_*^T K^{-1} f)\right). \quad (36)$$

References

- [1] AKCAY, S. Forecasting gamma-ray bursts using gravitational waves. [arXiv preprint arXiv:1808.10057](#) (2018).
- [2] DUVENAUD, D. The kernel cookbook: Advice on covariance functions. <https://www.cs.toronto.edu/~duvenaud/cookbook/>, 2014. Accessed: 2024-03-14.
- [3] HOY, C., AKÇAY, S., UILLIAM, J. M., AND THOMPSON, J. E. Incorporating multi-model uncertainty into gravitational-wave bayesian inference. [In preparation](#) (2024). To be published.
- [4] MAGGIORE, M. [Gravitational Waves. Volume 1: Theory and Experiments](#). Oxford University Press, Oxford, 2008.
- [5] OWEN, B. J. Search templates for gravitational waves from inspiraling binaries: Choice of template spacing. [arXiv preprint gr-qc/9511032](#) (1995). Submitted to Physical Review D: November 7, 1995.
- [6] RASMUSSEN, C. E., AND WILLIAMS, C. K. I. [Gaussian Processes for Machine Learning](#). MIT Press, Cambridge, MA, USA, 2006.