

On the Gaussian Process using JAX/NumPyro

HAJIME KAWAHARA ^{1,2}

¹Department of Space Astronomy and Astrophysics, ISAS/JAXA, 3-1-1, Yoshinodai, Sagami-hara, Kanagawa, 252-5210 Japan

²Department of Astronomy, Graduate School of Science, The University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, Japan

ABSTRACT

In this note, we consider the Gaussian process using JAX/NumPyro.

Keywords: JAX – NumPyro

1. GAUSSIAN PROCESS AS A STOCHASTIC MODEL

A Gaussian process is a stochastic process for a random variable \mathbf{x} that follows a multivariate normal distribution,

$$\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \Sigma), \quad (1)$$

where each component d_i is a time series. The covariance matrix of this multivariate normal distribution is defined as

$$\Sigma_{ij} = K_{ij}(a, \tau) = ak(|t_i - t_j|; \tau) \quad (2)$$

a function of the absolute difference between t_i and t_j . This formulation yields a Gaussian process with a correlation length τ . Various kernel functions $k(t, \tau)$ can be used. For example, the RBF kernel is given by

$$k_{\text{RBF}}(t; \tau) = \exp\left(-\frac{t^2}{2\tau^2}\right), \quad (3)$$

and the Matérn 3/2 kernel is defined as

$$k_{\text{M3/2}}(t; \tau) = \left(1 + \frac{\sqrt{3}t}{\tau}\right) e^{-\sqrt{3}t/\tau}. \quad (4)$$

Let us implement these kernels and generate some example data. The kernels can be coded in Python as follows:

```
1 import numpy as np
2 def RBF(t, tau):
3     Dt = t - np.array([t]).T
4     K = np.exp(-(Dt)**2 / 2 / (tau**2))
5     return K
6
7 def Matern32(t, tau):
8     Dt = t - np.array([t]).T
9     fac = np.sqrt(3.0) * np.abs(Dt) / tau
10    K = (1.0 + fac) * np.exp(-fac)
11    return K
```

19 Next, we use the ‘scipy.stats.multivariate_normal’ module to sample from a multivariate normal distribution. The
 20 sampling code is as follows:

```

1 from scipy.stats import multivariate_normal as smn
2 np.random.seed(seed=1) # Set random seed
3 N = 101
4 t = np.linspace(0, 10, N)
5 ave = np.zeros(N)
6 tau = 0.4
7 a = 1.0
8 cov = a * RBF(t, tau) # Alternatively, use Matern32(t, tau)
9 di = smn(mean=ave, cov=cov, allow_singular=True).rvs(1).T

```

21 To simulate observational noise, we add Gaussian noise with zero mean and a standard deviation σ :

```

1 sigma = 0.6
2 d = di + np.random.normal(0.0, sigma, len(di))

```

22 The generated data, shown in Figure 1, includes the original data without observational noise (solid line).

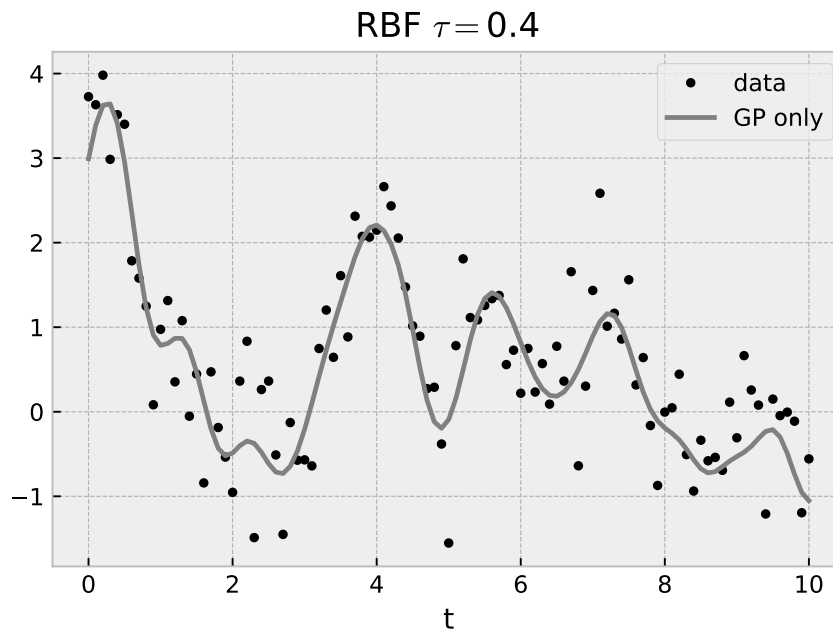


Figure 1. The solid line represents data without observational noise.

23 Now, we estimate the parameters τ , a , and σ from this data. With Gaussian noise added to the observations, the
 24 probabilistic model is expressed as:

$$\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \Sigma') \quad (5)$$

$$\Sigma' = K(a, \tau) + \sigma^2 I \quad (6)$$

27 The following implementation computes Σ' :

```

1 def modelcov(t, tau, a, sigma):
2     Dt = t - jnp.array([t]).T

```

```

3 K = a * jnp.exp(-(Dt)**2 / 2 / (tau**2)) + jnp.eye(N) * sigma**2
4 return K

```

28 The numpyro model for the Gaussian process using the modelcov function is defined as follows:

```

1 import jax.numpy as jnp
2 import numpyro
3 import numpyro.distributions as dist
4
5 def model(t, y):
6     sigma = numpyro.sample('sigma', dist.Exponential(1.))
7     tau = numpyro.sample('tau', dist.Exponential(1.))
8     a = numpyro.sample('a', dist.Exponential(1.))
9     cov = modelcov(t, tau, a, sigma)
10    numpyro.sample('y', dist.MultivariateNormal(loc=jnp.zeros(N),
11    covariance_matrix=cov), obs=y)

```

29 Here, the prior distributions for τ and σ are set to exponential distributions. To perform Hamiltonian Monte Carlo
30 (HMC) sampling using NUTS in numpyro, we proceed as follows:

```

1 from numpyro.infer import MCMC, NUTS
2 from jax import random
3 rng_key = random.PRNGKey(0)
4 rng_key, rng_key_ = random.split(rng_key)
5 num_warmup, num_samples = 1000, 2000
6 # Run NUTS.
7 kernel = NUTS(model)
8 mcmc = MCMC(kernel, num_warmup, num_samples)
9 mcmc.run(rng_key_, t=t, y=d)

```

31 The posterior distributions can be visualized using arviz, as shown below:

```

1 import arviz
2 refs = {}; refs["sigma"] = sigma; refs["tau"] = tau; refs["a"] = a
3 arviz.plot_pair(arviz.from_numpyro(mcmc), kind='kde',
4     divergences=False, marginals=True, reference_values=refs,
5     reference_values_kwargs={'color':'red', "marker": "o", "markersize":12})

```

32 The resulting visualization of the posterior distributions is shown in Figure 2.

33 1.1. Predictive Distribution and Credible Intervals

34 The model defined in Equation (5) assumes a mean of $\mathbf{x} = \mathbf{0}$. This can be further elucidated by computing the
35 credible intervals. Using the predictive posterior samples, the Highest Posterior Density Interval (HPDI) can be
36 computed as follows:

```

1 from numpyro.infer import Predictive
2 from numpyro.diagnostics import hpdi
3 # Extract posterior samples
4 posterior_a = mcmc.get_samples()['a']
5 posterior_tau = mcmc.get_samples()['tau']
6 posterior_sigma = mcmc.get_samples()['sigma']

```

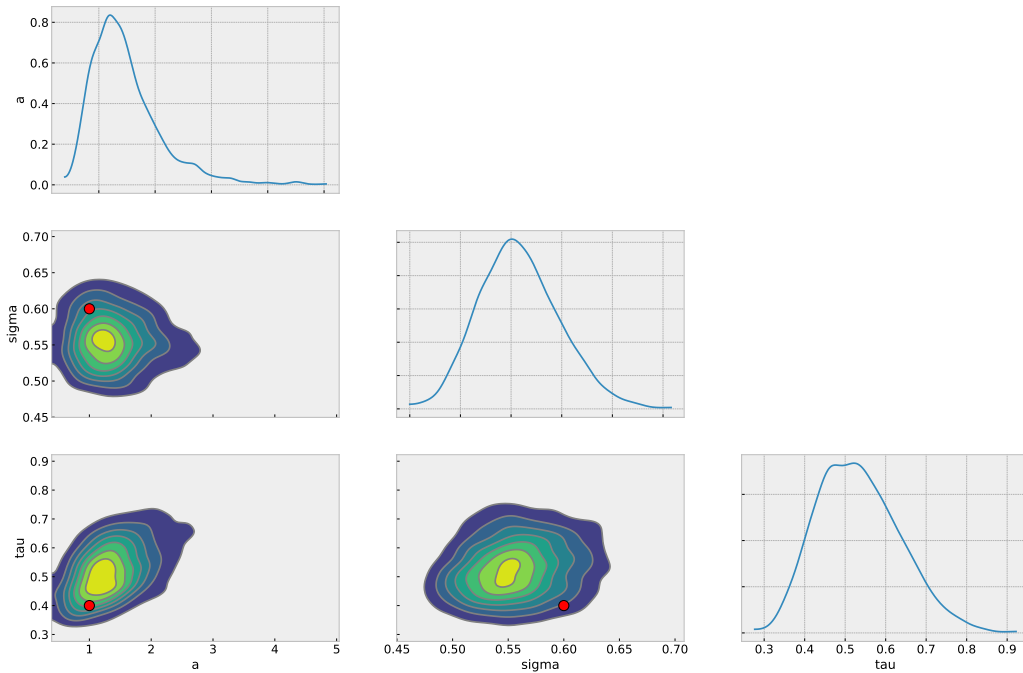


Figure 2.

```

7 # Predictive model
8 pred = Predictive(model, {'a': posterior_a, 'tau': posterior_tau,
9   'sigma': posterior_sigma}, return_sites=["y"])
10 predictions = pred(rng_key_, t=t, y=None)
11 mean_muy = jnp.mean(predictions["y"], axis=0)
12 hpdi_muy = hpdi(predictions["y"], 0.9)

```

The predicted mean and HPDI can then be plotted as follows:

```

1 fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(6, 3))
2 ax.plot(t, d, ".", color="black")
3 ax.plot(t, mean_muy, color="C0")
4 ax.fill_between(t, hpdi_muy[0], hpdi_muy[1], alpha=0.3, interpolate=True, color="C0")
5 plt.xlabel("t")
6 plt.ylabel("y")

```

The resulting plot, shown in Figure 3, provides an intuitive visualization of the data and the credible intervals.

1.2. Gaussian Processes as Priors for Model Parameters

The Gaussian process described above can also be used as a prior distribution for model parameters \mathbf{m} :

$$p(\mathbf{m}) = \mathcal{N}(\mathbf{0}, \Sigma) \quad (7)$$

Observations d_i can be expressed as:

$$d_i = m_i + \epsilon \quad (8)$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2) \quad (9)$$

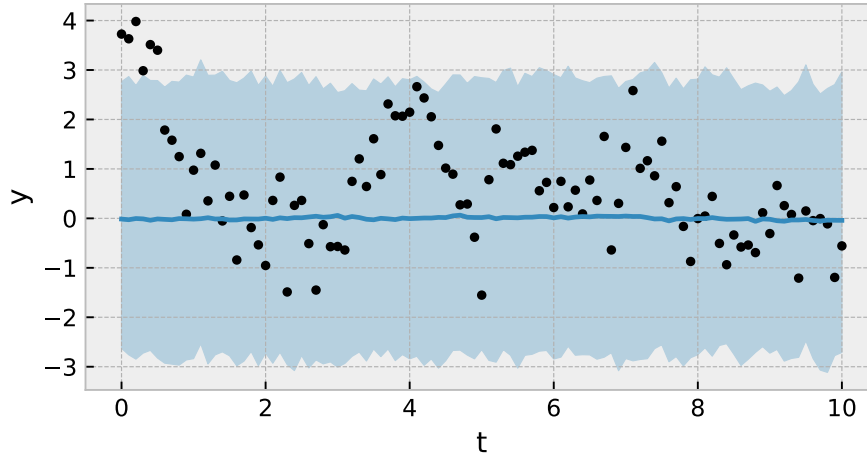


Figure 3.

45 where ϵ corresponds to observational noise. Alternatively, using the identity transformation for the model \mathbf{g} , this can
 46 be written as:

$$47 \quad \mathbf{g}(\mathbf{m}) = \mathbf{m} \quad (10)$$

$$48 \quad \mathbf{d} = \mathbf{g}(\mathbf{m}) + \epsilon \quad (11)$$

$$49 \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I) \quad (12)$$

50 The likelihood function becomes:

$$51 \quad p(\mathbf{d}|\mathbf{m}) = \mathcal{N}(\mathbf{d} - \mathbf{g}(\mathbf{m}), \sigma^2 I) = \mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I) \quad (13)$$

52 The covariance matrix Σ is parameterized as:

$$53 \quad \Sigma_{ij} = ak(|t_i - t_j|; \tau) \quad (14)$$

54 where a and τ are hyperparameters. These hyperparameters are assigned prior distributions (hyperpriors). While the
 55 posterior distribution of \mathbf{m} can be sampled directly using `numpyro`, it is also known that the posterior distribution can
 56 be derived analytically when the hyperparameters are fixed.

57 The exponent of a multivariate normal distribution can be expressed as:

$$58 \quad -2 \log \mathcal{N}(\mathbf{m}|\boldsymbol{\mu}, \Sigma) = (\mathbf{m} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{m} - \boldsymbol{\mu}) = \mathbf{m}^\top \Sigma^{-1} \mathbf{m} - 2\mathbf{m}^\top \Sigma^{-1} \boldsymbol{\mu} + \text{const.} \quad (15)$$

59 If the density $p(\mathbf{m})$ is represented as:

$$60 \quad -2 \log p(\mathbf{m}) = \mathbf{m}^\top P \mathbf{m} - 2\mathbf{m}^\top \mathbf{q} + \text{const}, \quad (16)$$

61 then comparing this to the multivariate normal distribution yields:

$$62 \quad p(\mathbf{m}) = \mathcal{N}(\mathbf{m}|P^{-1}\mathbf{q}, P^{-1}). \quad (17)$$

63 Now, the posterior distribution is expressed as:

$$64 \quad p(\mathbf{m}|\mathbf{d}) \propto p(\mathbf{d}|\mathbf{m})p(\mathbf{m}) = \mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I)\mathcal{N}(\mathbf{0}, \Sigma) \quad (18)$$

and by expanding the exponent terms with respect to \mathbf{m} , we can write:

$$-2 \log [\mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I) \mathcal{N}(\mathbf{0}, \Sigma)] = \mathbf{m}^\top (\Sigma^{-1} + \sigma^{-2} I) \mathbf{m} - 2\sigma^{-2} \mathbf{m}^\top \mathbf{d} + \text{const.} \quad (19)$$

Additionally, using the relationship:

$$(\Sigma^{-1} + \sigma^{-2} I)^{-1} = \Sigma(I + \sigma^{-2} \Sigma)^{-1} \quad (20)$$

we find:

$$p(\mathbf{m}|\mathbf{d}) = \mathcal{N}(\Sigma(\sigma^2 I + \Sigma)^{-1} \mathbf{d}, \Sigma(I + \sigma^{-2} \Sigma)^{-1}) \quad (21)$$

1.3. Hyperparameter Sampling Using HMC-NUTS

The earlier sampling of τ and σ via HMC-NUTS can now be interpreted in this framework. These parameters are hyperparameters, denoted as $\boldsymbol{\theta} = (a, \tau, \sigma)^\top$. Incorporating $\boldsymbol{\theta}$ into the probabilistic model, the marginalized likelihood over \mathbf{m} becomes:

$$p(\mathbf{d}|\boldsymbol{\theta}) = \frac{p(\mathbf{d}|\mathbf{m}, \boldsymbol{\theta}) p(\mathbf{m}, \boldsymbol{\theta})}{p(\mathbf{m}|\mathbf{d}, \boldsymbol{\theta})} \quad (22)$$

By focusing on the terms involving \mathbf{d} in the exponent:

$$-2 \log p(\mathbf{d}|\boldsymbol{\theta}) = -2 \log p(\mathbf{d}|\mathbf{m}, \boldsymbol{\theta}) + 2 \log p(\mathbf{m}|\mathbf{d}, \boldsymbol{\theta}) + \text{const} \quad (23)$$

$$= -2 \log \mathcal{N}(\mathbf{d} - \mathbf{m}; \sigma^2 I) + 2 \log \mathcal{N}(\Sigma(\sigma^2 I + \Sigma)^{-1} \mathbf{d}, \Sigma(I + \sigma^{-2} \Sigma)^{-1}) \quad (24)$$

$$= \mathbf{d}^\top (\sigma^2 I + \Sigma)^{-1} \mathbf{d} + \text{const.} \quad (25)$$

we obtain:

$$p(\mathbf{d}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \Sigma + \sigma^2 I) \quad (26)$$

This matches Equation (5) when $\Sigma = K(\tau)$. Thus, by providing a prior $p(\boldsymbol{\theta})$, the marginalized posterior distribution:

$$p(\boldsymbol{\theta}|\mathbf{d}) \propto p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (27)$$

is sampled via HMC-NUTS. Figure 3 shows the marginalized posterior distribution of the hyperparameters τ and σ . That is, we sampled

$$\boldsymbol{\theta}_k^\dagger \sim p(\boldsymbol{\theta}|\mathbf{d}) \propto p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta}). \quad (28)$$

From the joint posterior distribution:

$$p(\mathbf{m}, \boldsymbol{\theta}|\mathbf{d}) = p(\mathbf{m}|\boldsymbol{\theta}, \mathbf{d}) p(\boldsymbol{\theta}|\mathbf{d}) \quad (29)$$

and the posterior distribution of \mathbf{m} given $\boldsymbol{\theta}_k^\dagger$:

$$p(\mathbf{m}|\boldsymbol{\theta}_k^\dagger, \mathbf{d}) = \mathcal{N}(\boldsymbol{\mu}_k, K_k) \quad (30)$$

$$\boldsymbol{\mu}_k = K(a_k^\dagger, \tau_k^\dagger) ((\sigma_k^\dagger)^2 I + K(a_k^\dagger, \tau_k^\dagger))^{-1} \mathbf{d} \quad (31)$$

$$K_k = K(a_k^\dagger, \tau_k^\dagger) (I + (\sigma_k^\dagger)^{-2} K(a_k^\dagger, \tau_k^\dagger))^{-1} \quad (32)$$

we sample $(\mathbf{m}_k^\dagger, \boldsymbol{\theta}_k^\dagger)$ from $p(\mathbf{m}, \boldsymbol{\theta}|\mathbf{d})$. The implementation for $\boldsymbol{\mu}_k$ and K_k is as follows:

```

1 def muGP(tau, a, sigma):
2     cov = a * RBF(t, tau)
3     IKw = sigma**2 * np.eye(N) + cov
4     A = scipy.linalg.solve(IKw, d, assume_a="pos")
5     return cov @ A
6
7 def covGP(tau, a, sigma):
8     cov = a * RBF(t, tau)
9     IKw = np.eye(N) + cov / sigma**2
10    IKw = scipy.linalg.inv(IKw)
11    return cov @ IKw

```

94 Sampling \mathbf{m} using the posterior samples of $\boldsymbol{\theta}$ is done as follows:

```

1 import tqdm
2 Ns = len(posterior_sigma)
3 np.random.seed(seed=1)
4 marr = []
5 for i in tqdm.tqdm(range(0, Ns)):
6     sigmas = float(posterior_sigma[i])
7     taus = float(posterior_tau[i])
8     a_s = float(posterior_a[i])
9     ave = muGP(taus, a_s, sigmas)
10    cov = covGP(taus, a_s, sigmas)
11    mk = smn(mean=ave, cov=cov, allow_singular=True).rvs(1).T
12    marr.append(mk)
13 marr = np.array(marr)

```

95 The credible interval for \mathbf{m} is then computed as follows:

```

1 mean_muy = np.mean(marr, axis=0)
2 hpdi_muy = hpdi(marr, 0.9)

```

96 Plotting the results, as shown in Figure 4, reveals the 90% credible interval for the model parameter \mathbf{m} . Note that
97 this interval is for \mathbf{m} , the model parameter, and does not include observational noise:

```

1 fig = plt.figure(figsize=(6, 3))
2 ax = fig.add_subplot(111)
3 ax.plot(t, d, ".", color="black")
4 ax.plot(t, mean_muy, color="C1")
5 ax.fill_between(t, hpdi_muy[0], hpdi_muy[1], alpha=0.3, interpolate=True, color="C1")
6 plt.xlabel("t")
7 plt.ylabel("y")

```

98 1.4. Predictions at New Data Points

99 To predict at $t = t^*$, the posterior distribution of \mathbf{m}^* is:

$$100 \quad p(\mathbf{m}^* | \mathbf{m}) = \mathcal{N}(K_{\times}^{\top} K^{-1} \mathbf{m}, K_{*} - K_{\times}^{\top} K^{-1} K_{\times}) \quad (33)$$

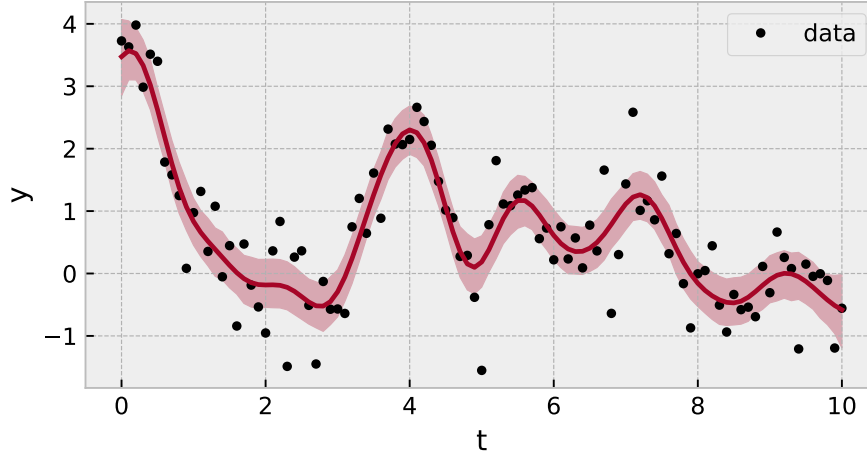


Figure 4.

101 where

$$K_{ij} = ak(|t_i - t_j|; \tau) \quad (34)$$

$$102 (K_{\times})_{ij} = ak(|t_i - t_j^*|; \tau) \quad (35)$$

$$103 (K_{*})_{ij} = ak(|t_i^* - t_j^*|; \tau). \quad (36)$$

104 Similarly,

$$105 p(\mathbf{m}^* | \mathbf{d}) = \mathcal{N}(K_{\times}^{\top} K_{\sigma}^{-1} \mathbf{d}, K_{*} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}) \quad (37)$$

106 where

$$107 (K_{\sigma})_{ij} = ak(|t_i - t_j|; \tau) + \sigma^2 \delta_{i,j} \quad (38)$$

108 If observational noise is included, as

$$109 \mathbf{d}^* = \mathbf{m}^* + \boldsymbol{\epsilon}, \quad (39)$$

110 the predictive distribution becomes:

$$111 p(\mathbf{d}^* | \mathbf{d}) = \mathcal{N}(K_{\times}^{\top} K_{\sigma}^{-1} \mathbf{d}, K_{*,\sigma} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}) \quad (40)$$

112 where:

$$113 (K_{*,\sigma})_{ij} = ak(|t_i^* - t_j^*|; \tau) + \sigma^2 \delta_{i,j} \quad (41)$$

114 Sampling from the predictive distributions can be implemented as follows:

```

1 def mucovGPx(t, td, tau, a, sigma):
2     cov = a * RBF(t, tau) + sigma**2 * np.eye(N)
3     covx = a * RBFx(t, td, tau)
4     covxx = a * RBF(td, tau)
5     # For predictions including observational noise
6     covxx += sigma**2 * np.eye(len(td))
7     IKw = cov
8     A = scipy.linalg.solve(IKw, d, assume_a="pos")
9     IKw = scipy.linalg.inv(IKw)
10    return covx @ A, covxx - covx @ IKw @ covx.T # Mean and covariance

```


116 This enables predictions and credible intervals to be computed, as shown in Figure 5, with darker regions representing
 117 credible intervals for \mathbf{m}^* and lighter regions for \mathbf{d}^* .

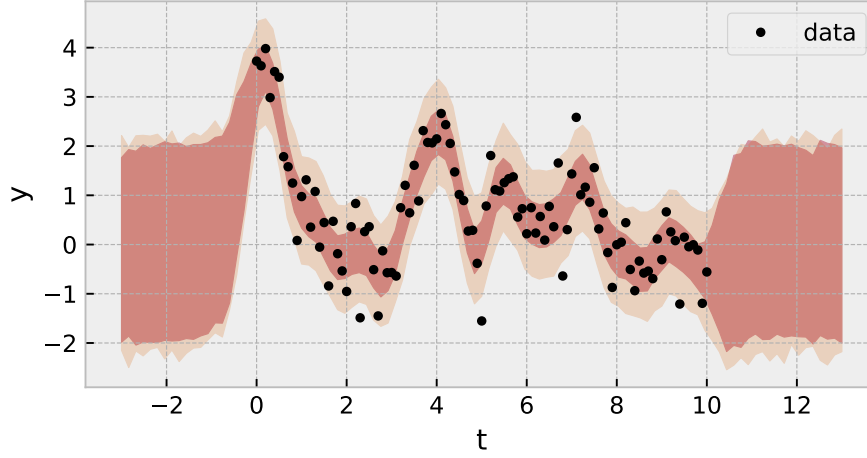


Figure 5.

118 2. MODEL FITTING WITH GAUSSIAN PROCESS NOISE

119 So far, we have considered a model generated by a Gaussian process with a zero mean:

$$120 \quad \mathbf{m} \sim \mathcal{N}(\mathbf{0}, \Sigma(t)). \quad (42)$$

121 Here, we explore a case where the data is generated by a Gaussian process with a mean function dependent on t .
 122 Specifically, we consider:

$$123 \quad \mathbf{m} \sim \mathcal{N}(\mathbf{f}(t), \Sigma(t)), \quad (43)$$

124 where $\mathbf{f}(t)$ is a vectorized version of $f(t)$, applied to each element of $\mathbf{t} = (t_0, t_1, \dots, t_{N-1})$. For $f(t)$, let us assume
 125 the following functional form:

$$126 \quad f(t) = ke^{-(t-T_0)^2/2s^2} \sin(2\pi t/P). \quad (44)$$

127 While we write this as $f(t)$ for simplicity, we may also use the notation $f(t; \boldsymbol{\theta})$ to explicitly represent its dependence
 128 on the parameters $\boldsymbol{\theta} = (T_0, k, s, P)$.

129 Using the RBF kernel for the covariance matrix in Equation (43), we can generate data as shown in Figure 6. Here,
 130 $\tau = 3$, and we observe a relatively smooth trend overlaying $f(t)$.

131 Fitting such a model using HMC allows us to model noise correlated with the signal $f(t)$, including observational
 132 noise, and estimate the parameters of $f(t)$.

133 If $\mathbf{f}(t)$ is known, the predictions, both without and with observational noise, are given by:

$$134 \quad p(\mathbf{m}^* | \mathbf{d}) = \mathcal{N}(\mathbf{f}(t^*) + K_{\times}^{\top} K_{\sigma}^{-1} (\mathbf{d} - \mathbf{f}(t)), K_{*} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}), \quad (45)$$

$$135 \quad p(\mathbf{d}^* | \mathbf{d}) = \mathcal{N}(\mathbf{f}(t^*) + K_{\times}^{\top} K_{\sigma}^{-1} (\mathbf{d} - \mathbf{f}(t)), K_{*,\sigma} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}). \quad (46)$$

136 The parameters $\boldsymbol{\theta} = (T_0, k, s, P)$ of $\mathbf{f}(t)$ are sampled using HMC. For example, in the latter case, predictions can be
 137 obtained by sampling each $\boldsymbol{\theta}_k^{\dagger}$ as follows:

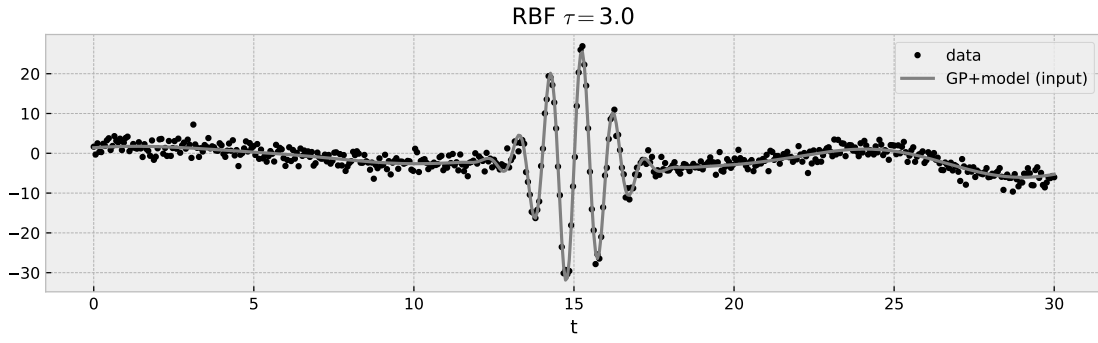


Figure 6. The solid line represents data without observational noise.

$$\mathbf{d}_k^* \sim \mathcal{N}(\mathbf{f}(t^*; \boldsymbol{\theta}_k^\dagger) + K_\times^\top K_\sigma^{-1}(\mathbf{d} - \mathbf{f}(t; \boldsymbol{\theta}_k^\dagger)), K_{*,\sigma} - K_\times^\top K_\sigma^{-1} K_\times). \quad (47)$$

Thus, we can perform Gaussian process fitting that incorporates the model $f(t)$, as shown in Figure 7.

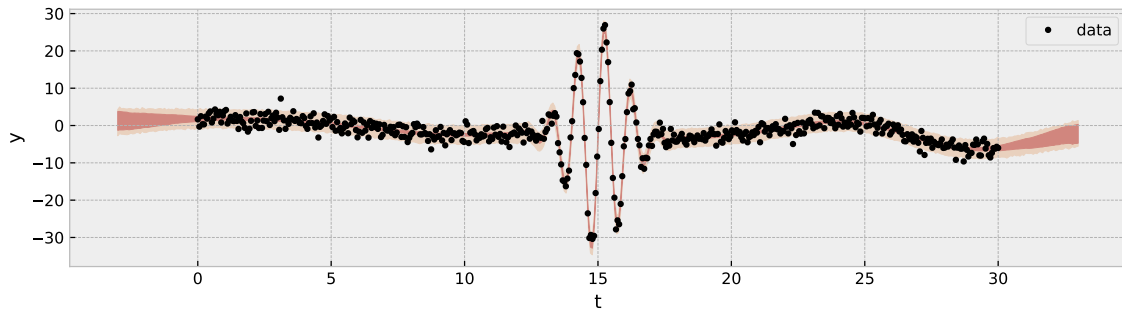


Figure 7.

Ref: Appendix in [Kawahara et al. \(2022\)](#)

REFERENCES

- ¹⁴¹ Kawahara, H., Kawashima, Y., Masuda, K., et al. 2022,
¹⁴² ApJS, 258, 31, doi: [10.3847/1538-4365/ac3b4d](https://doi.org/10.3847/1538-4365/ac3b4d)