

# Part III Astrostatistics: Example Sheet 3 Solution

## Example Class: Friday, 13 Mar 2020, 12:00pm, MR5

### 1 Doubly Lensed Quasar Time Delay Estimation

Quasar light curves (brightness time series)  $f(t)$  (in magnitudes) are often modelled using the Ornstein-Uhlenbeck (O-U) process, which is described mathematically by a stochastic differential equation of the form:

$$df(t) = \tau^{-1}[c - f(t)] dt + \sigma dW_t \quad (1)$$

where the second term is a Brownian motion (continuous-time limit of a random walk), with the variability scaled by  $\sigma$ , and the first term is a drag that tends to return the brightness back to the mean level  $c$ . The O-U process is a Gaussian process

$$f(t) \sim \mathcal{GP}(m(t), k(t, t')) \quad (2)$$

with mean function  $m(t) = c$  and covariance function or kernel:

$$\text{Cov}[f(t), f(t')] = k(t, t') = A^2 \exp(-|t - t'|/\tau) \quad (3)$$

with characteristic amplitude  $A^2 = \tau\sigma^2/2$ . The characteristic timescale for the quasar brightness to revert to the mean  $c$  is  $\tau$ . Hence, astronomers often call this a “damped random walk”. In a doubly-lensed quasar system, two images of the same quasar are observed. However, their brightness time series will have a time delay and magnification relative to each other due to the gravitational lensing effects. Find in the accompanying dataset, measurements of the brightness time series of two images of a lensed quasar,  $y_1(t)$  and  $y_2(t)$ . Assume the measurement errors are Gaussian with the given standard deviations. Where possible, write down and derive the relevant equations before you implement them in code.

1. Plot the data. For each image ( $y_1$  or  $y_2$ ) time series separately, fit an O-U process by optimising the marginal likelihood to estimate  $c$ ,  $A$ , and  $\tau$  for each image. Are these estimates consistent between the two time series? Estimate the overall relative magnification factor (difference in magnitudes) between the two images. (The relative multiplicative magnification  $\mu$  due to the gravitational lens is related to the magnitude shift by  $\Delta m = -2.5 \log_{10} \mu$ ).

**Solution: See Figures 1, 2, and code. Optimising the GP marginal likelihood (derived below) for  $y_1$ , we find**

$$c_1 = -0.0063; A_1 = 0.0911; \tau_1 = 41.5204.$$

**Optimising the GP marginal likelihood (derived below) for  $y_2$ , we find**

$$c_2 = 0.1053; A_2 = 0.1026; \tau_2 = 40.5140.$$

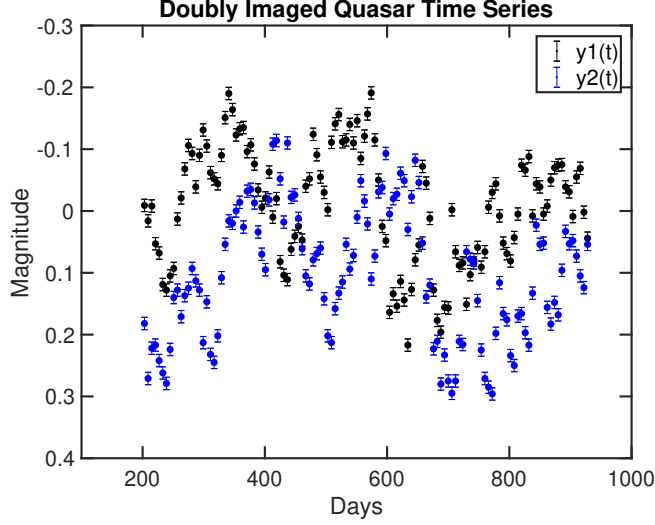


Figure 1: Plot of both raw time series datasets.

We find that  $A$  and  $\tau$  are generally consistent between the two time series. The mean levels are different by  $\Delta c = c_2 - c_1 \approx 0.10$ , suggesting that the magnitude shift is  $\Delta m \approx 0.10$ .

Let  $y_1$  be the vector of time series measurements of quasar image 1 at observation times  $t_o$ . That is,  $y_1^i$  is the measured magnitude of image 1 at time  $t_o^i$ . Let  $f_o$  be the vector of latent values of the underlying function at the observation times  $t_o$ . Since  $f(t)$  is a realisation of a Gaussian process,  $f_o$  has a multivariate Gaussian prior distribution,

$$f_o \sim N(\mathbf{1}c, K(t_o, t_o))$$

where  $\mathbf{1}$  is a vector of ones of the same dimension as  $f_o$ , and  $K(t_o, t_o)$  is a matrix, whose  $(i, j)$ th entry is an evaluation of the kernel  $k(t_o^i, t_o^j)$  between the pair of observation times  $t_o^i, t_o^j$ , and depends on  $A, \tau$ . The measurement covariance matrix  $W$  can be constructed by placing the stated measurement variances  $\sigma_i^2$  on the diagonal  $W_{ii}$ , and assuming the measurement covariances are zero (as they are not reported). Then the measurement process can be described by:

$$y_1 | f_o \sim N(f, W)$$

Using the properties of multivariate Gaussians, and/or the fact that  $y_1$  results from the sum of two Gaussian random vectors,  $f_o$  and  $\epsilon \sim N(0, W)$ , we have the sampling distribution

$$\begin{aligned} P(y_1 | c, A, \tau^2) &= \int P(y_1, f_o | c, A, \tau^2) df_o = \int N(y_1 | f_o, W) N(f_o | \mathbf{1}c, K(t_o, t_o)) df_o \\ &= N(y_1 | \mathbf{1}c, K(t_o, t_o) + W) = N(y_1 | \mathbf{1}c, \Sigma(A, \tau)) \end{aligned}$$

where we define  $\Sigma(A, \tau) \equiv K(t_o, t_o) + W$ . When viewed as a function of the parameters  $c, A, \tau$ , with the data fixed, this is the marginal likelihood function

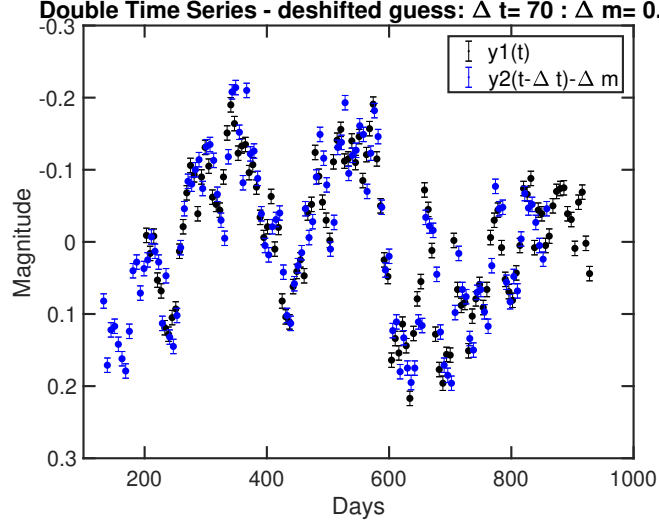


Figure 2: Plot of time series  $y_1$  and  $y_2$  shifted by guesses for  $\Delta t$  and  $\Delta m$ .

(because we have effectively marginalised out  $f_o$ ). The marginal likelihood for  $y_2$  can be derived the same way.

To optimise the marginal likelihood, it is best to numerically minimise the negative log likelihood function with respect to  $c, A, \tau$ :

$$-\log P(y_1 | c, A, \tau^2) = \frac{1}{2} \log \det \Sigma(A, \tau) + \frac{1}{2} (y_1 - 1c)^T \Sigma(A, \tau)^{-1} (y_1 - 1c)$$

It is inadvisable in the first term to directly compute  $\det \Sigma(A, \tau)$  and then take the log, because the determinant may likely be near zero. Instead, since  $\Sigma$  must be a proper covariance matrix, we can perform a Cholesky decomposition,  $LL^T = \Sigma$ , where  $L$  is the lower triangular Cholesky factor. Then we can use the identity,

$$\log \det \Sigma = 2 \sum_{i=1}^N \log L_{ii}$$

along the diagonal of the Cholesky factor to directly compute the log determinant. Now that we have computed the Cholesky factor, we can also use it to compute the quadratic form in the second term. Letting  $r = y_1 - 1c$ ,

$$r^T \Sigma^{-1} r = r^T [LL^T]^{-1} r = r^T [(L^{-1})^T L^{-1}] r = z^T z$$

where  $z$  is the solution to the linear system of equations  $Lz = r$ . While this is mathematically equivalent to  $z = L^{-1}r$ , it can be solved using forward substitution without inverting  $L$ . Thus, the quadratic form can be computed without directly inverting  $\Sigma$ , which is usually an expensive and often numerically unstable operation.

2. Fixing the values of  $c$ ,  $A$ , and  $\tau$  you found for each image separately, overplot random light curves drawn from the GP prior on each separate time series dataset. Use the Gaussian Process machinery to estimate the underlying light curve of each image separately. Plot the expectation and standard deviation of the posterior prediction as a function of time.

**Solution:** See Figures 3, 5, 4, 6, and the code.

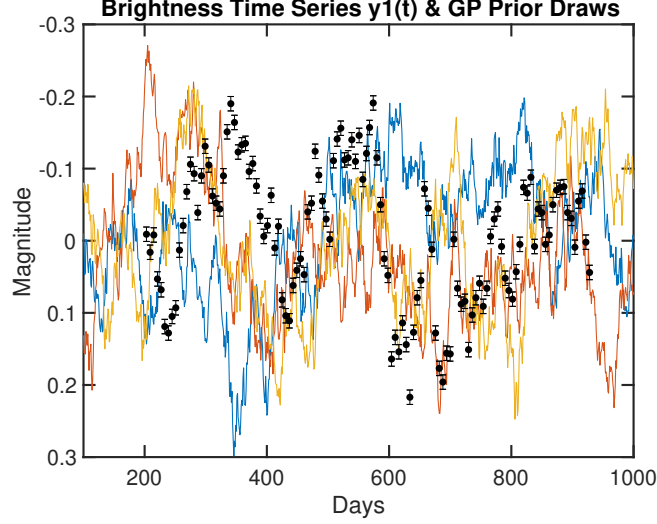


Figure 3: Plot of  $y_1$  data and draws from the GP prior with hyperparameters optimised for  $y_1$ .

Let  $t_g$  be a grid of times on which we want to infer the latent function  $f(t)$ . Let  $f_g$  be the values of the latent function at those times, i.e.  $f_g^i = f(t_g^i)$ . The joint distribution of  $y_1$  and  $f_g$  is

$$\begin{pmatrix} y_1 \\ f_g \end{pmatrix} \sim N \left( \begin{pmatrix} 1c \\ 1c \end{pmatrix}, \begin{pmatrix} K(t_o, t_o) + W & K(t_o, t_g) \\ K(t_g, t_o) & K(t_g, t_g) \end{pmatrix} \right).$$

The covariance matrix components can be derived by considering  $y_1 = f_o + \epsilon$ , where  $\epsilon \sim N(0, W)$ , and computing  $\text{Cov}[y_1, y_1]$ ,  $\text{Cov}[y_1, f_g]$ ,  $\text{Cov}[f_o, f_g]$ ,  $\text{Cov}[f_g, f_g]$ , etc. Then using conditional property of multivariate Gaussians random vectors, we can compute the posterior of  $f_g$  given  $y_1$ .

$$f_g | y_1 \sim N(\mathbb{E}[f_g | y_1], \text{Var}[f_g | y_1])$$

$$\mathbb{E}[f_g | y_1] = 1c + K(t_g, t_o) \Sigma^{-1} (f_g - 1c)$$

$$\text{Var}[f_g | y_1] = K(t_g, t_g) - K(t_g, t_o) \Sigma^{-1} K(t_o, t_g)$$

and similarly for  $f_g | y_2$ .

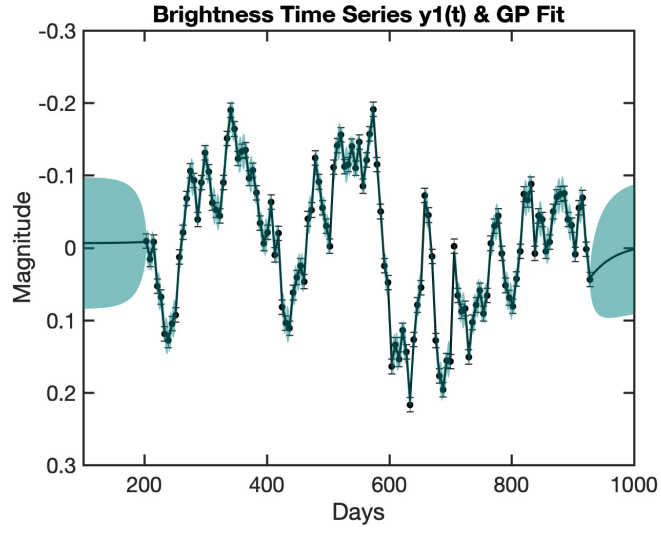


Figure 4: Plot of  $y_1$  data and the GP posterior fit to  $y_1$ .

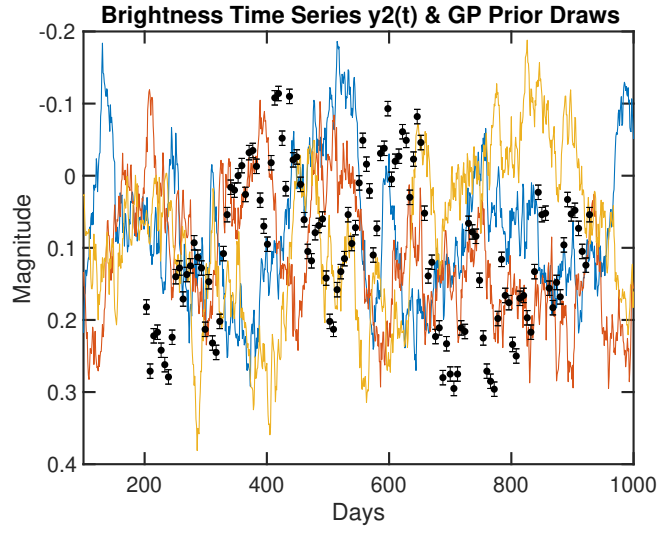


Figure 5: Plot of  $y_2$  data and draws from the GP prior with hyperparameters optimised for  $y_2$ .

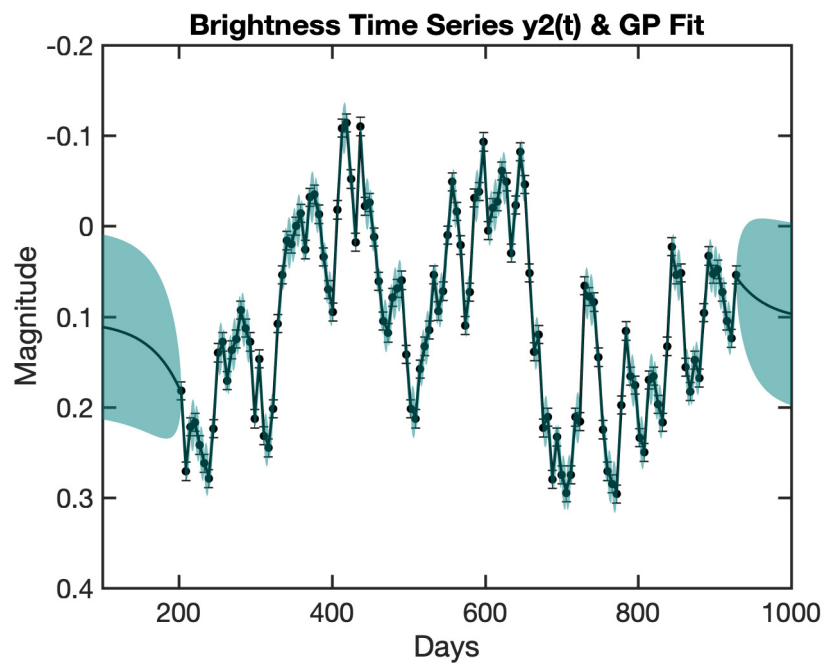


Figure 6: Plot of  $y_2$  data and the GP posterior fit to  $y_2$ .

3. Derive a likelihood function for the two time series considered jointly, as two copies of the same realisation of the GP but shifted in time by the time delay  $\Delta t$ , and the magnification factor  $\Delta m$  (both relative to  $y_1$ ), and measured with noise at the observed times. Thus  $y_1(t)$  is a noisy measurement of  $f(t)$  and  $y_2(t)$  is a noisy measurement of  $f(t - \Delta t) + \Delta m$ . Using suitable non-informative priors, write down a posterior density  $P(\Delta t, \Delta m, c, A, \tau | \mathbf{y}_1, \mathbf{y}_2)$ .

**Solution:** Suppose we knew the true time delay  $\Delta t$  and magnitude shift  $\Delta m$ . Then we could construct the combined vectors:

$$\tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 - \Delta m \end{pmatrix}$$

$$\tilde{\mathbf{t}} = \begin{pmatrix} \mathbf{t}_o \\ \mathbf{t}_o - \Delta t \end{pmatrix}$$

such that  $\tilde{\mathbf{y}}$  contains noisy measurements of the same latent function  $\tilde{f}$  at times  $\tilde{\mathbf{t}}$  (i.e.  $\tilde{f}^i = f(\tilde{t}^i)$ ).

$$\tilde{f} \sim N(\mathbf{1}c, \mathbf{K}(\tilde{\mathbf{t}}, \tilde{\mathbf{t}}))$$

$$\tilde{\mathbf{y}}_1 | \tilde{f} \sim N(\tilde{f}, \tilde{\mathbf{W}})$$

where

$$\tilde{\mathbf{W}} = \begin{pmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{W} \end{pmatrix}.$$

Therefore the marginal likelihood function is

$$P(\tilde{\mathbf{y}} | \Delta t, \Delta m, c, A, \tau) = N(\tilde{\mathbf{y}} | \mathbf{1}c, \mathbf{K}(\tilde{\mathbf{t}}, \tilde{\mathbf{t}}) + \tilde{\mathbf{W}})$$

which again should be computed in log-space according to the advice above.

4. Estimate  $\Delta t$  and  $\Delta m$ . Beware that the log likelihood is highly multimodal, so it is important to find the major mode. You may fix the O-U process parameters to reasonable values found previously, or estimate them jointly with  $(\Delta t, \Delta m)$ .

**Solution:** First, I optimise the previous marginal likelihood for both time series, using as initial guesses  $c = c_1$ ,  $\Delta m = c_2 - c_1$ ,  $A = \frac{1}{2}(A_1 + A_2)$ , and  $\tau = \frac{1}{2}(\tau_1 + \tau_2)$ . After optimising the marginal likelihood and using the observed Fisher information to derive parameter uncertainties, I find

$$\widehat{\Delta t} = 74.96 \pm 0.64$$

$$\widehat{\Delta m} = 0.10 \pm 0.29$$

$$\hat{c} = 0.01 \pm 0.04$$

$$\hat{A} = 0.10 \pm 0.01$$

$$\hat{\tau} = 43.87 \pm 0.19.$$

A contour plot of the marginal likelihood near the optimum is shown in Fig. 7. To explore the multi-modality of the marginal likelihood more closely, in Fig. 8, I have plotted the log of the profile likelihood:

$$L_{\text{prof}}(\Delta t) = \arg \max_{\Delta m, c, A, \tau} P(\tilde{\mathbf{y}} | \Delta t, \Delta m, c, A, \tau).$$

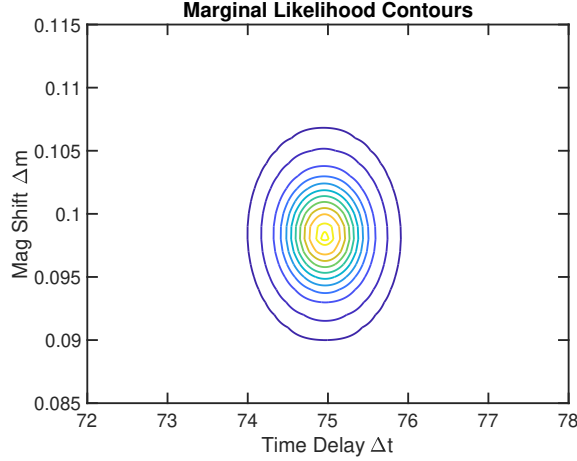


Figure 7: Plot of the marginal likelihood of  $(\Delta t, \Delta m)$  with the other hyperparameters fixed to their MLE values.

In the code, I implemented a Metropolis-within-Gibbs algorithm to sample the posterior of the shifts with the hyperparameters fixed to their MLE values. I use flat improper priors on the time delay and magnitude shift.

$$P(\Delta t, \Delta m | \tilde{\mathbf{y}}, \hat{c}, \hat{A}, \hat{\tau}) \propto P(\tilde{\mathbf{y}} | \Delta t, \Delta m, \hat{c}, \hat{A}, \hat{\tau})$$

I ran 4 chains in parallel for  $10^4$  iterations, each initialised with the MLEs plus some random jitter. The two proposal jump scales in each parameter was tuned to give 30 – 50% acceptance rates. The Gelman-Rubin ratios for both parameters was  $< 1.0004$ . In Fig. 9, I show the sample autocorrelation function of the slowest mixing parameter ( $\Delta m$ ). The estimated autocorrelation time was 18 iterations, which I take as a thinning factor. I remove 20% of each chain as burn-in and thin the remainder by a factor of 18 (keeping only every 18th sample). Then I concatenated the chains to compute posterior inferences on the parameters. Figure 10, shows a scatter plot of the joint distribution of  $(\Delta t, \Delta m)$  samples, showing that the parameters are well-constrained and uncorrelated in the posterior. Figure 11 show the 1D marginal posteriors distributions as 1D histograms of each parameter estimates. The posterior means and standard deviations are shown.



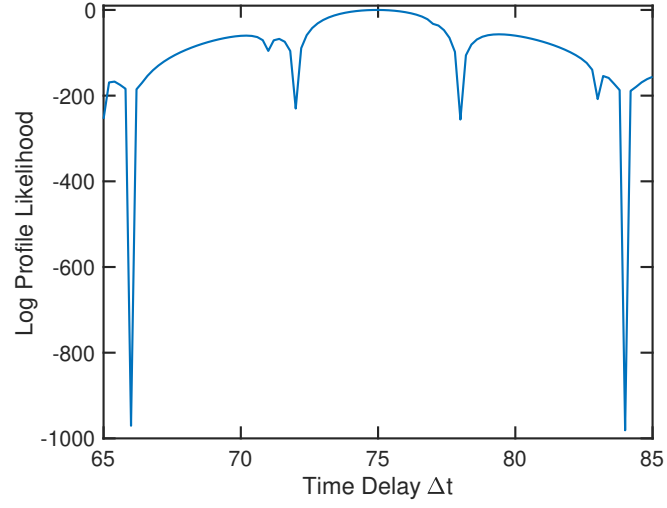


Figure 8: **Plot of the profile likelihood of  $\Delta t$ .**

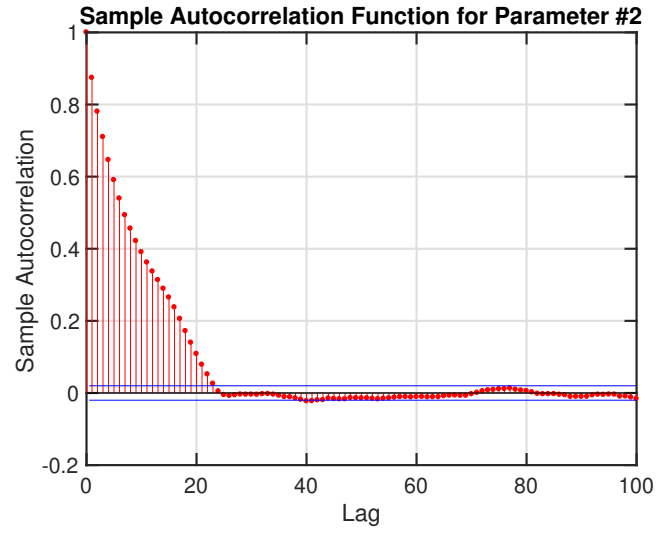


Figure 9: **Plot of the MCMC autocorrelation function of  $\Delta m$ .**

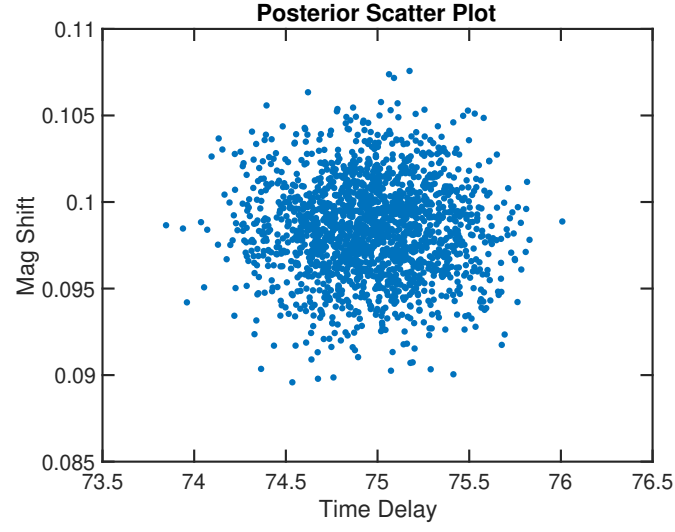


Figure 10: **Plot of the posterior distribution of  $(\Delta t, \Delta m)$  MCMC samples.**

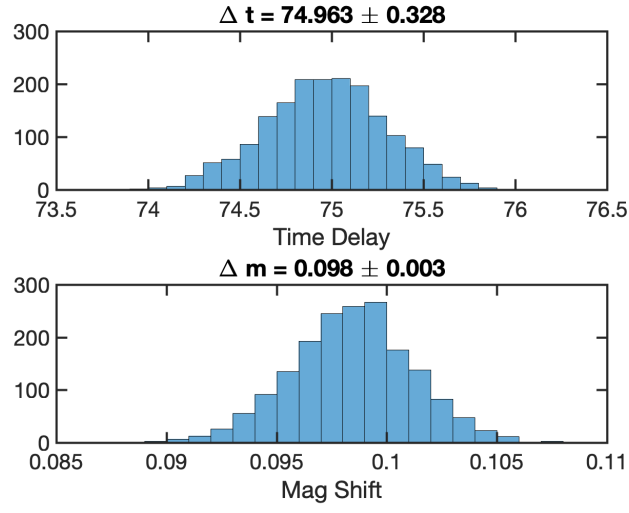


Figure 11: **Plot of the marginal posterior histograms of  $(\Delta t, \Delta m)$  MCMC samples.**

5. Overplot the two time series datasets, with  $y_2$  shifted to the  $y_1$  frame by subtracting the estimated  $\Delta t$  and  $\Delta m$ . Now using the O-U parameters you found, plot the posterior estimate of the underlying light curve using the two combined datasets.

**Solution:** Figure 12 shows the posterior GP fit to both time series using the estimated time delay and magnitude shift.

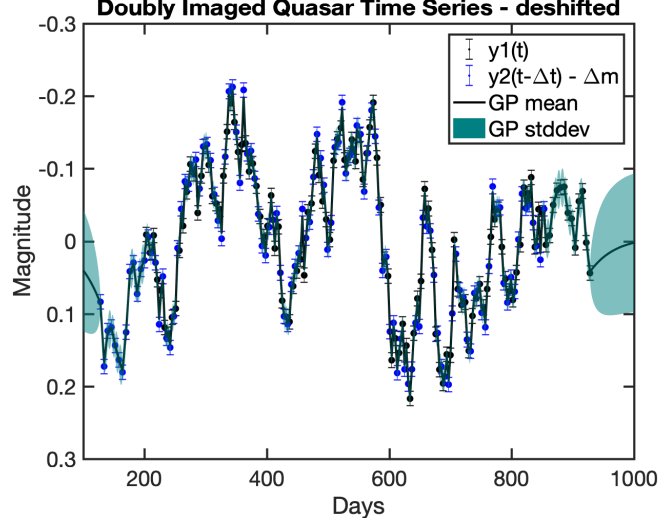


Figure 12: Plot of the marginal posterior histograms of  $(\Delta t, \Delta m)$  MCMC samples.

With  $c, A, \tau, \Delta t, \Delta m$  set to their fitted estimates, we can compute the posterior of the latent function  $f_g$  evaluated on grid points  $t_g$ :

$$\begin{aligned} \begin{pmatrix} \tilde{y} \\ f_g \end{pmatrix} &\sim N \left( \begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} K(\tilde{t}, \tilde{t}) + W & K(\tilde{t}, t_g) \\ K(t_g, \tilde{t}) & K(t_g, t_g) \end{bmatrix} \right) \\ f_g | \tilde{y} &\sim N(\mathbb{E}[f_g | \tilde{y}], \text{Var}[f_g | \tilde{y}]) \\ \mathbb{E}[f_g | \tilde{y}] &= 1c + K(t_g, \tilde{t})[K(\tilde{t}, \tilde{t}) + W]^{-1}(f_g - 1c) \\ \text{Var}[f_g | \tilde{y}] &= K(t_g, t_g) - K(t_g, \tilde{t})[K(\tilde{t}, \tilde{t}) + W]^{-1}K(\tilde{t}, t_g). \end{aligned}$$

*Numerical Clue:* A proper covariance matrix admits a Cholesky decomposition:  $\Sigma = LL^T$ , where  $L$  is the lower triangular Cholesky factor. The log of the determinant  $|\Sigma| = \det \Sigma$  can stably be computed from Equation A.18 of Rasmussen & Williams. If you have computed the Cholesky factor  $L$ , solutions  $\mathbf{x}$  to linear equations of the form  $\Sigma \mathbf{x} = \mathbf{b}$ , i.e.  $\mathbf{x} = \Sigma^{-1} \mathbf{b}$ , can be stably computed using forward/backward substitution, rather than by directly inverting  $\Sigma$ , as described in Rasmussen & Williams, §A.4.