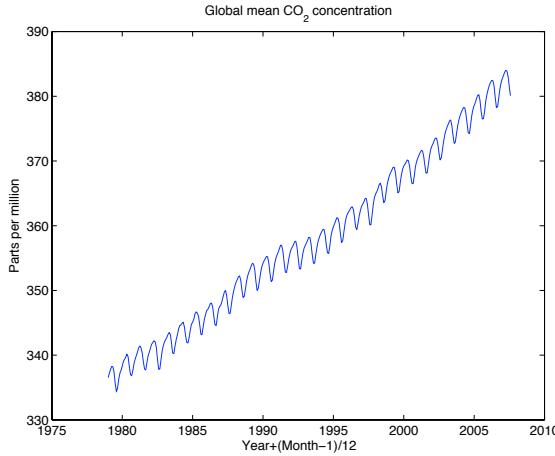


Bayesian linear and Gaussian process regression. The following time series of monthly mean global CO₂ concentrations can be obtained from the file `co2.txt` (original data obtained from <http://www.esrl.noaa.gov/gmd/ccgg/trends>):



We will apply Bayesian linear and Gaussian process regression to predict the CO₂ concentration $f(t)$ as a function of time t , where $t = \text{Year} + (\text{Month} - 1)/12$.

- (a) First we model the function using linear regression, that is, using the functional form

$$f(t) = at + b + \epsilon(t),$$

with i.i.d. noise residual $\epsilon(t) \sim \mathcal{N}(0, 1)$ and prior $a \sim \mathcal{N}(0, 10^2)$, $b \sim \mathcal{N}(360, 100^2)$. Compute (using MATLAB or another package) the posterior mean and covariance over a and b given the CO₂ data.

- (b) Let $a_{\text{MAP}}, b_{\text{MAP}}$ be the MAP estimate in the question above. The residual is the difference between the observed function values and the predicted mean function values

$$g_{\text{obs}}(t) = f_{\text{obs}}(t) - (a_{\text{MAP}}t + b_{\text{MAP}}),$$

where $f_{\text{obs}}(t)$ is the observed value of the CO₂ concentration at time t .

Plot $g_{\text{obs}}(t)$. Do you think these residuals conform to our prior over $\epsilon(t)$? State, with justifications, which characteristics of the residual you think do or do not conform to our prior belief.

- (c) Write a MATLAB (or other language) function to generate samples drawn from a GP. Specifically, given a covariance kernel function $k(\cdot, \cdot)$ and a vector of input points \mathbf{x} , return a function $f(\mathbf{x})$ evaluated on the input points \mathbf{x} drawn randomly from a GP with the given covariance kernel and with zero mean.
- (d) Test your function by plotting sample functions drawn from the following kernel, for various settings of the hyperparameters

$$k(s, t) = \theta^2 \left(\exp \left(-\frac{2 \sin^2(\pi(s-t)/\tau)}{\sigma^2} \right) + \phi^2 \exp \left(-\frac{(s-t)^2}{2\eta^2} \right) \right) + \zeta^2 \delta_{s=t} \quad (1)$$

Describe the characteristics of the drawn functions, and how the characteristics of the functions depend on the parameters. [5 marks]

covariance_mtx = \sum

$$\Sigma = \text{diag} = \text{chol} \left(\mathbb{I} + \begin{pmatrix} 0.001 & & \\ & \ddots & 0 \\ 0 & & 0.001 \end{pmatrix} \right).$$

$$\text{Sample} = \Sigma \times \begin{pmatrix} \text{v}(0,1) \\ \vdots \\ \text{v}(n,n) \end{pmatrix} + \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} \text{v}(0,1) \\ \vdots \\ \text{v}(n,n) \end{pmatrix}$$

- (e) Suppose we were to consider modelling the residual function $g(t)$ using a zero mean GP with the covariance kernel above. Based on the plot of $g(t)$ and your explorations in the preceding part, what do you think will be suitable values for the hyperparameters of k ?
- (f) [Bonus] Extrapolate the CO₂ concentration levels to 2020 using the GP with covariance kernel k of eqn 1, and your chosen parameter values. Specifically, compute the predictive mean and variance of the residual $g(t)$ for every month between September 2007 and December 2020 given the observed residuals $g_{\text{obs}}(t)$. Plot the means and one standard deviation error bars of the extrapolated CO₂ concentration levels

$$f(t) = a_{\text{MAP}} t + b_{\text{MAP}} + g(t)$$

along with the observed CO₂ levels. Does the behaviour of the extrapolation conform to your expectations? How sensitive are your conclusions to settings of the kernel hyperparameters?

- (g) [Bonus] Why is the above procedure not fully Bayesian? How would we go about modelling $f(t)$ in a Bayesian framework?

(a) Consider a general Bayesian Regression problem :

$$y = \phi^T w + \varepsilon,$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$w \sim \mathcal{N}(m_0, S_0).$$

i.e. we have a prior that $P(w) = \mathcal{N}(m_0, S_0)$.

By the Baye's theorem, the posterior over parameters is :

$$P(w|X, y) = \frac{P(y|X, w) \cdot P(w)}{P(y|X)} \quad \dots \textcircled{*}$$

independent of w

where $P(y|X) = \int P(y|X, w) \cdot P(w) dw$. is the evidence
(marginal likelihood).

$$P(y|X, w) = \mathcal{N}(y|\phi w, \sigma^2 I).$$

CLAIM : posterior of parameters : $P(w|X, y) = \mathcal{N}(w|S_N, m_N)$,

$$\text{where } S_N = (S_0^{-1} + \sigma^{-2} \phi^T \phi)^{-1}$$

$$m_N = S_N (S_0^{-1} m_0 + \sigma^{-2} \phi^T y).$$

Proof : From $\textcircled{*}$, $P(w|X, y) \propto P(y|X, w) \cdot P(w)$

Take log of RHS :

$$\begin{aligned}
 & \log \mathcal{N}(y|\phi w, \sigma^2 I) + \log \mathcal{N}(w|m_0, S_0). \quad \text{Independent of } w \\
 &= -\frac{1}{2} \left[\sigma^{-2} (y - \phi w)^T (y - \phi w) + (w - m_0)^T S_0^{-1} (w - m_0) \right] + \text{const} \\
 &= -\frac{1}{2} \left[\cancel{\sigma^{-2} y^T y} - \cancel{2\sigma^{-2} y^T \phi w} - \cancel{\sigma^{-2} w^T \phi^T \phi w} + \cancel{w^T S_0^{-1} w} - \cancel{2m_0^T S_0^{-1} w} + \cancel{(m_0^T S_0^{-1} m_0)} \right] + \text{const} \\
 &= -\frac{1}{2} \left[\underbrace{w^T (\sigma^{-2} \phi^T \phi + S_0^{-1}) w}_{\text{const}} - \cancel{2(\sigma^{-2} y^T \phi + m_0^T S_0^{-1}) w} \right] + \text{const}' \\
 &= -\frac{1}{2} \left[\underbrace{w^T (\sigma^{-2} \phi^T \phi + S_0^{-1}) w}_{\text{const}} - \cancel{2(\sigma^{-2} \phi^T y + S_0^{-1} m_0)^T w} \right] + \text{const}'
 \end{aligned}$$

$$\text{Hence } P(w|X,y) \propto \exp \left\{ -\frac{1}{2} \left[\underbrace{w^\top (G^{-2} \phi^\top \phi + S_0^{-1}) w}_{\text{green}} - \underbrace{2(G^{-2} \phi^\top y + S_0^{-1} m_0)^\top w}_{\text{red}} \right] \right\}$$

$$\begin{aligned} \text{Compare } P(w|X,y) &\propto \exp \left\{ -\frac{1}{2} (w - m_N)^\top S_N^{-1} (w - m_N) \right\} \\ &= \exp \left\{ -\frac{1}{2} \left[\underbrace{w^\top S_N^{-1} w}_{\text{green}} - \underbrace{2m_N^\top S_N^{-1} w}_{\text{red}} + m_N^\top S_N^{-1} m_N \right] \right\} \end{aligned}$$

$$\Rightarrow \begin{cases} S_N^{-1} = (G^{-2} \phi^\top \phi + S_0^{-1}) \\ m_N^\top S_N^{-1} = (G^{-2} \phi^\top y + S_0^{-1} m_0)^\top \end{cases} \Leftrightarrow \begin{cases} S_N = (G^{-2} \phi^\top \phi + S_0^{-1})^{-1} \\ m_N = S_N (G^{-2} \phi^\top y + S_0^{-1} m_0) \end{cases} \quad \square$$

- - - - - Proof completed - - - - - - - -

In this case,

$$G^2 = 1, \quad S_0 = \begin{pmatrix} 10^2 & 0 \\ 0 & 100^2 \end{pmatrix}, \quad m_0 = \begin{pmatrix} 0 \\ 360 \end{pmatrix}, \quad \phi = \begin{pmatrix} \dots & \ell & \dots \\ \dots & 11 & \dots \end{pmatrix}$$

Plug in matlab, we get m_N (posterior-mean)
 S_N (posterior-covariance)

posterior_covariance =

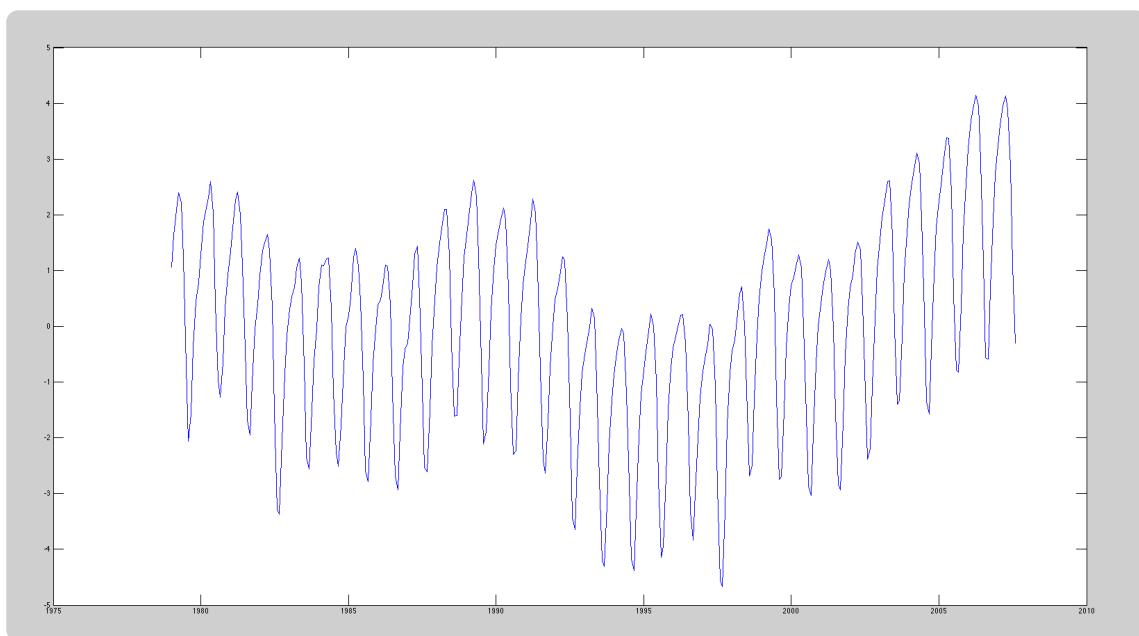
$$\begin{aligned} 1.0e+02 * \\ 0.000000417453539 &\quad -0.000832106417939 \\ -0.000832106417939 &\quad 1.658659376255700 \end{aligned}$$

posterior_mean =

$$\begin{aligned} 1.0e+03 * \\ 0.001572033536338 \\ -2.775586453638971 \end{aligned}$$

>>

(b)



The residuals DO NOT confirm our prior over $\varepsilon(t)$.

It is given $\varepsilon(t) \sim N(0, 1)$.

But the variance we get for ε is 3.324, which is much greater than 1.

**g_obs_variance =
3.323762842551880
fx >>**

Also this can be spotted from the plot.

(c) Function 1 : covarianceKernelFunction

```
covarianceKernelFunction.m x SampleFromGP.m x generateCovarianceMatrix.m x +  
1 function [ cov ] = covarianceKernelFunction(X1,X2,parameters)  
2  
3     % split the parameters  
4     theta = parameters(1);  
5     tau = parameters(2);  
6     sigma = parameters(3);  
7     phi = parameters(4);  
8     eta = parameters(5);  
9     zeta = parameters(6);  
10  
11    if X1 == X2  
12        tail = zeta^2 * 1;  
13    else  
14        tail = 0;  
15    end  
16  
17    cov = theta^2 * (exp((-2) * (sin(pi*(X1-X2)/tau)^2))/(sigma^2) ...  
18        + phi^2 * exp((- (X1-X2)^2)/(2*eta^2))) + tail;  
19  
20 end
```

Function 2 : generateCovarianceMatrix

```
covarianceKernelFunction.m x SampleFromGP.m x generateCovarianceMatrix.m x +  
1 function covariance_matrix = generateCovarianceMatrix(t,parameters)  
2     covariance_matrix = zeros(length(t),length(t));  
3     for i = 1:length(t)  
4         for j = 1:length(t)  
5             covariance_matrix(i,j) = covarianceKernelFunction(t(i),t(j),parameters);  
6         end  
7     end
```

Function 3 : SampleFromGP

```
covarianceKernelFunction.m x SampleFromGP.m x generateCovarianceMatrix.m x +  
1 function sample = SampleFromGP(mean,covariance_matrix)  
2  
3     % Cholesky decomposition  
4     sig = chol(covariance_matrix + eye(size(covariance_matrix))*0.00000001, 'lower');  
5     % generate a vector (seed) with mean 0, variance 1  
6     unit_vector = randn(size(covariance_matrix,1),1);  
7     % project to generate the final sample  
8     sample = sig * unit_vector + mean;  
9  
10 end
```

(d) θ : determines the whole amplitude

large $\theta \Rightarrow$ large oscillation amplitude.

σ : determines length of the "wiggles" in the function.

τ : determines the repetition of the peak.

i.e. the "small" period. In this case, $\tau = 1$ as it represents each year as a period.

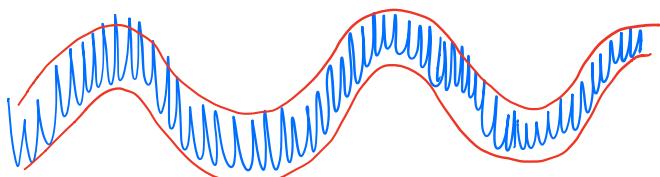
ϕ : the magnitude of how "twisted" the function is.

γ : the period of twist. "large" period.

ξ : magnitude of noise.

Consider an example:

τ, σ determines the "small" oscillation in BLUE

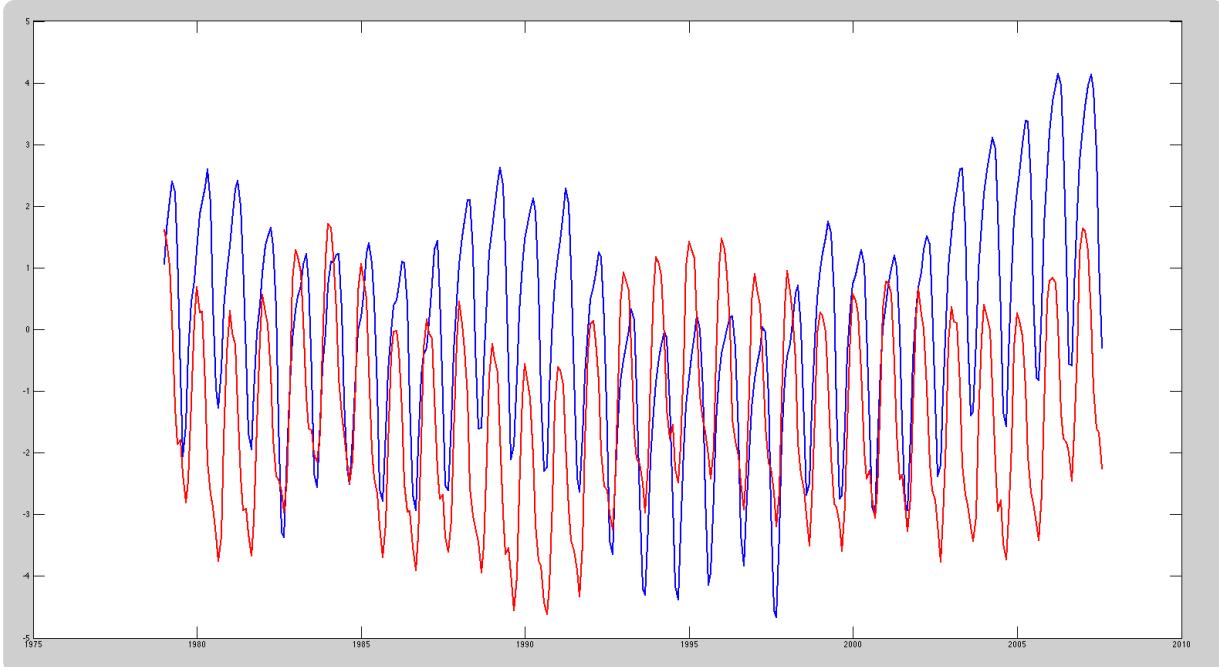


ϕ, γ determine the "large" oscillation in RED

```
(e) 19 - theta = 1.85;
20 - tau = 1;
21 - sigma = 1.3;
22 - phi = 0.52;
23 - eta = 1.36;
24 - zeta = 0.08;
25 - parameters = [theta tau sigma phi eta zeta];
```

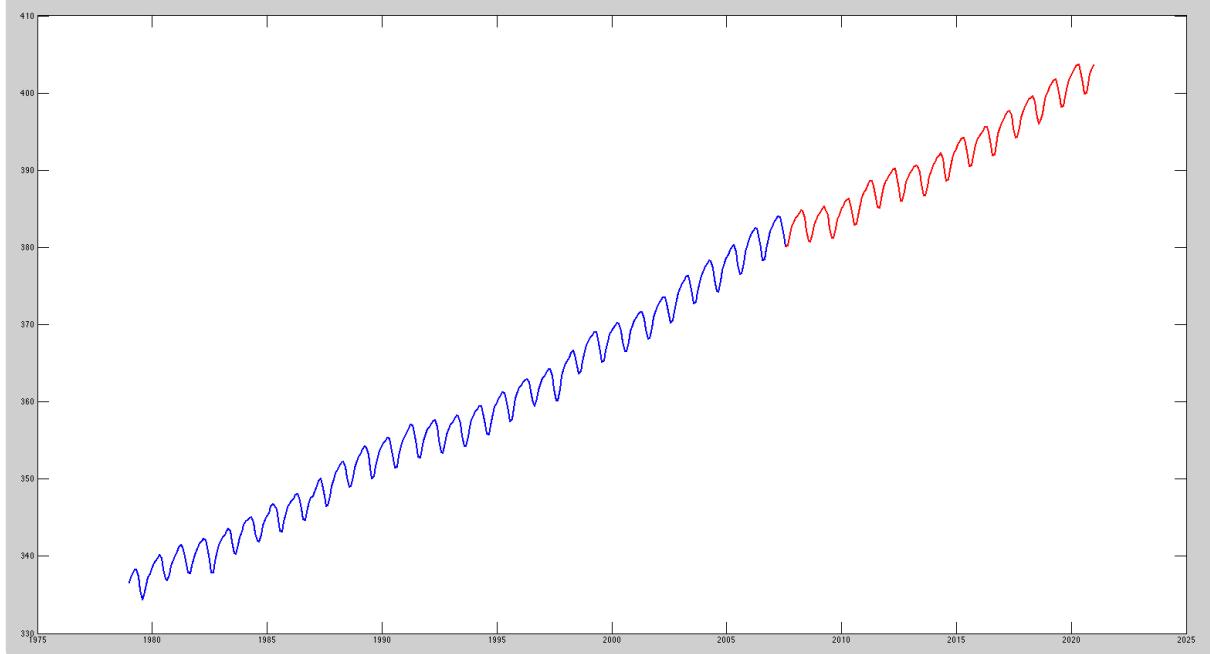
My choice : $\Theta = 1.85$, $\underline{\tau = 1}$, $\sigma = 1.3$, $\phi = 0.52$, $\eta = 1.36$
 \downarrow
 one year

$$\xi = 0.08$$



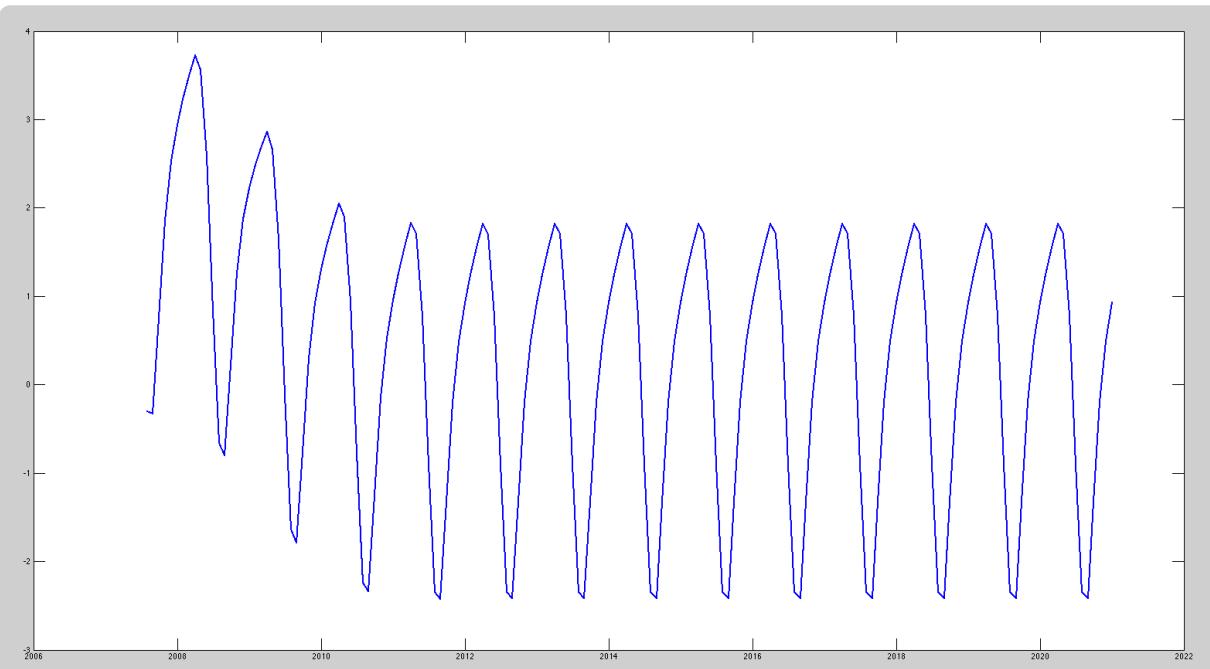
Plot modelled error with true error with the given hyperparameter.
 Red Blue

(f)

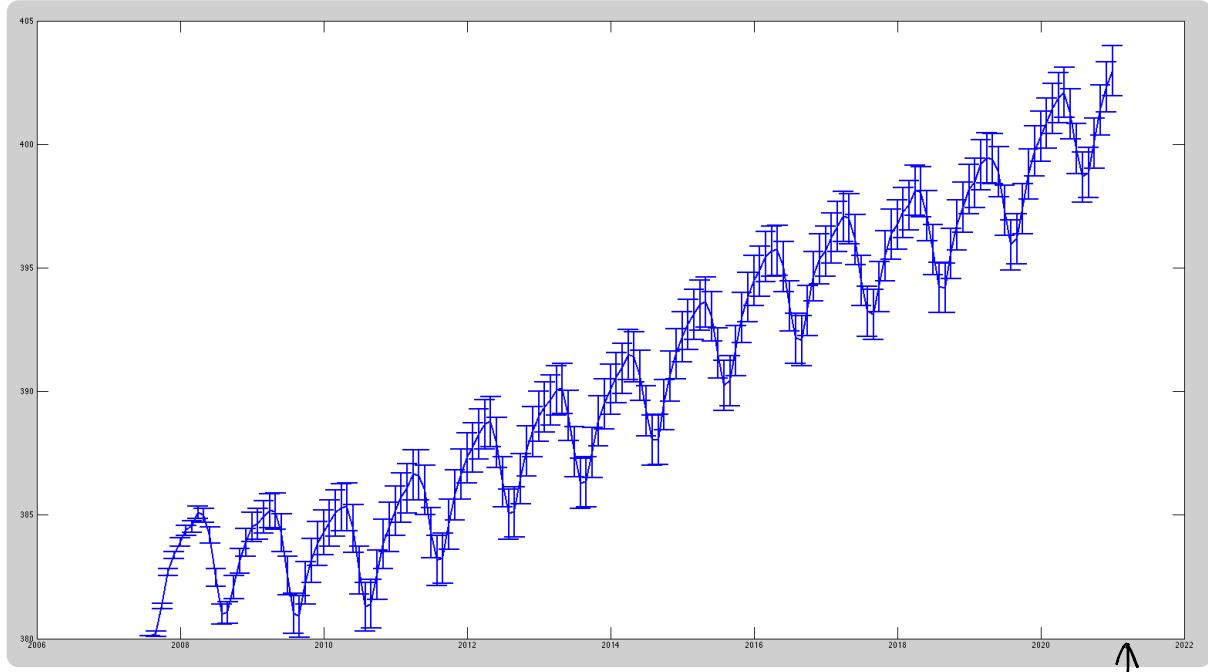


BLUE : original data Red : predicted data

Predictive mean plot



With error bar \checkmark (errors are diagonal of covariance matrix)



Narrow error bar at beginning ↑

Wide error bar, further from data ↑

The extrapolation conforms my expectation.

* it oscillates every year in the "small" period.

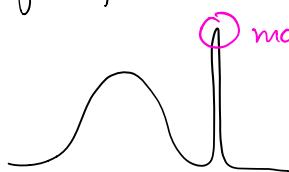
* there is a "large" period, which slightly twists the curve.

the "large period" is around 18 years, which has been shown by the prediction.

How sensitive: It is very sensitive to the parameters.

if we significantly change the parameters,
the shape will be changed significantly.

(g). The procedure is NOT fully Bayesian since MAP is a point estimate, it might find the highest peak not the largest density.



How to modify to make it fully Bayesian?

Starting with a non-linear regression problem $f(x) = \omega^T \phi(x)$.

$$Y^T | X \sim N(0_n, \tau^2 \phi^T \phi + \sigma^2 I_n)$$

Predict a single new test data by:

$$y | \underline{x}, \underline{y}, \underline{\underline{x}} \sim N(\tilde{K}_{\underline{x}\underline{x}} \tilde{K}_{\underline{x}\underline{x}}^{-1} \underline{y}^T, \tilde{K}_{\underline{x}\underline{x}} - \tilde{K}_{\underline{x}\underline{x}} \tilde{K}_{\underline{x}\underline{x}}^{-1} \tilde{K}_{\underline{x}\underline{x}})$$

and choose a specific kernel $k(x_1, x_2) = \phi(x_1)^T \phi(x_2)$.

Q(3) Appendix: the MATLAB script

```

1 % load data
2 - clear, clc, close all, load co2.txt -ascii; data = co2; format long
3 - t = data(:,1) + (data(:,2) - 1) / 12;
4 - f_obs = data(:,3); X = [t'; ones(1,length(t))];
5 - w_covariance = [10^2 0; 0 100^2];
6 - w_mean = [0 360]'; sigma = 1; % variance of error
7 - posterior_covariance = (w_covariance^(-1) + sigma^(-2) * X * X')^(-1);
8 - posterior_mean = posterior_covariance * (w_covariance^(-1) * w_mean + sigma^(-2) * X * f_obs );
9 - posterior_covariance
10 - posterior_mean
11 - a_map = posterior_mean(1); b_map = posterior_mean(2);
12 - figure();
13 - g_obs = f_obs - (a_map * t + b_map);
14 - g_obs_variance = std(g_obs)^2
15 - plot(t, g_obs); hold on
16 % hyperparameters:
17 - theta = 1.85; tau = 1; sigma = 1.3; phi = 0.52; eta = 1.36; zeta = 0.08;
18 - parameters = [theta tau sigma phi eta zeta];
19 - covariance_matrix = generateCovarianceMatrix(t,parameters);
20 - mean = zeros(length(t),1); % mean is zero, given in the question
21 - error = SampleFromGP(mean,covariance_matrix);
22 - % sample = mvnrnd(mean,covariance_matrix);
23 - plot(t,error,'r')
24 - % generate date to 2020
25 - t_future = (t(length(t)):1/12:2021)';
26 - t_full = [t;t_future];
27 - n = length(t); n_future = length(t_future); n_full = length(t_full);
28 - covariance_matrix_full = generateCovarianceMatrix(t_full,parameters);
29 - % split the full covariance matrix to use the conditional property:
30 - Kaa = covariance_matrix;
31 - Kab = covariance_matrix_full(1:n,n+1:n_full);
32 - Kba = covariance_matrix_full(n+1:n_full,1:n);
33 - Kbb = covariance_matrix_full(n+1:n_full,n+1:n_full);
34 - % distribution of the predicted data:
35 - prediction_mean = Kba * Kaa^(-1) * g_obs;
36 - prediction_covariance = Kbb - Kba * Kaa^(-1) * Kab;
37 - prediction_error = SampleFromGP(prediction_mean, prediction_covariance);
38 - prediction_data = a_map * t_future + ones(n_future,1) * b_map + prediction_error;
39 - figure()
40 - plot(t,data(:,3)); hold on
41 - plot(t_future,prediction_data,'r'); hold on
42 - figure()
43 - % plot errorbars
44 - plot(t_future,prediction_data,'r');
45 - err = sqrt(diag(prediction_covariance)); % diagonal of the prediction covariance matrix
46 - errorbar(t_future,prediction_data, err)

```