

**CS146 Final Projects
Viet Hoang Tran Duong
Fall 2019**

Executive Summary: CO₂ measurement modeling

Topic: CO₂ measurements of the world has been increasing steadily since the start of the industrial revolution in the 18th century. We attempt to build a statistical model to predict future CO₂ measurements to better prepare for future outcomes.

Data: Mauna Loa CO₂ measurements from ice core measurements.

Attempts: Create a statistical model to explain the data and forecast how the measurements will look like in 40 years, predict when the CO₂ level will be considered high risk (>450 ppm).

The outline of the projects is as followed:

1. Data preprocessing	3
2. Model 1: Example model and its flaws	4
3. Model 2: Updated version of model 1	6
4. Model 3: New proposed model.	7
5. Model comparison	9
6. Predictions for the year 2058: the first measurement in 2058	10
7. Predictions regarding the high-risk CO ₂ level (≥ 450 ppm)	11
8. Shortcomings	12
a. The process designing model	
i. Assumptions	12
ii. Priors	13
iii. Parameters	13
iv. Observed vs. Unobserved	13
b. Test statistics	13
9. Discussion & Recommendations	17
10. HCs Applications	18
11. References	18

Part 1: Data preprocessing

Step 1: Subtract the date data from the first date (1958-03-29) to have the Δt : time difference from the initial date. We do this because the model takes inputs as numeric values, not as date values. Hence, transforming the data would make it more convenient for running pystan.

Step 2: Visualize the data.

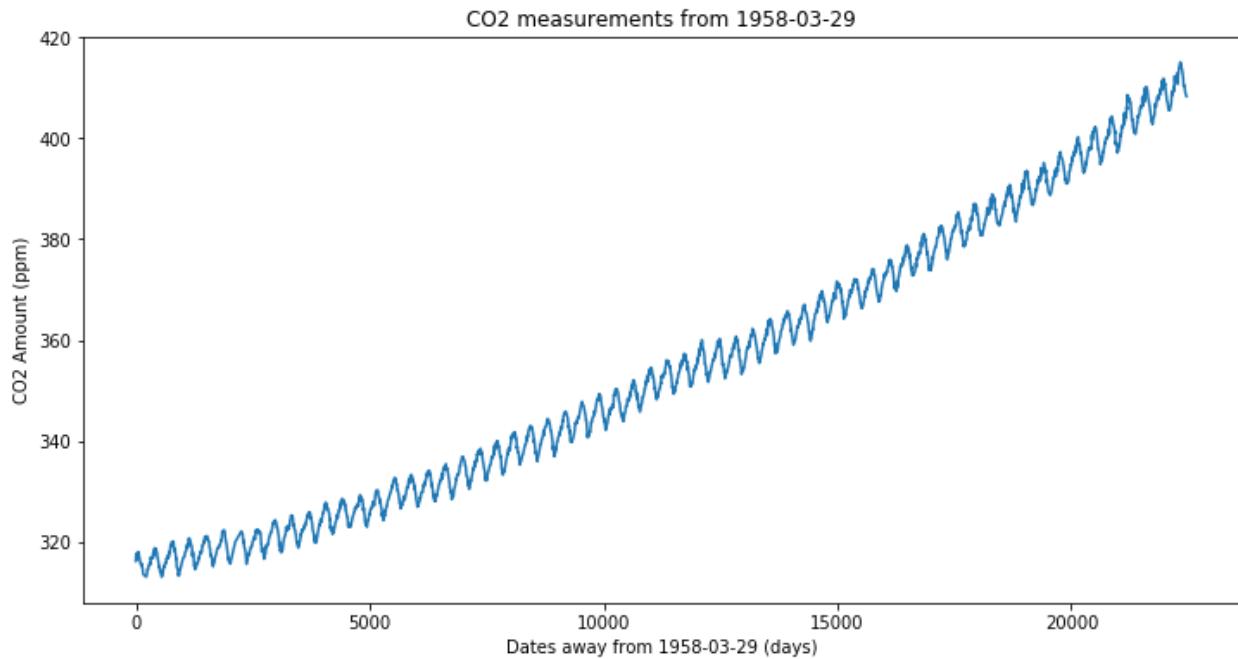


Fig. 1: CO₂ measurements from 1958-03-29

Step 3: Data partitioning

I subset the data into the first 70% data (2197 values) as the input for computing posterior distributions over the parameters using Stan, and the remaining latter 30% (942 values) to validate our prediction of future functions values.

Part 2: Model 1: The example model and its flaws

$$p(x_t|\theta) = N\left(c_0 + c_1 t + c_2 \cos\left(\frac{2\pi t}{365.25} + c_3\right), c_4^2\right)$$

- First, without any constraints on the parameters c_0, c_1, c_2, c_3, c_4 , the model diverges, with very correlated samples (based on the n_eff and Rhat values)

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
c_0	310.31	1.4e-3	0.09	310.14	310.25	310.31	310.37	310.48	3593	1.0
c_1	3.7e-3	1.6e-7	9.4e-6	3.7e-3	3.7e-3	3.7e-3	3.7e-3	3.7e-3	3658	1.0
c_2	2.71	0.09	0.14	2.47	2.57	2.72	2.84	2.94	2	2.54
c_3	2.92	2.07	2.92	6.6e-5	9.0e-4	2.9	5.85	5.88	2	218.23
c_4	2.0	0.06	0.09	1.87	1.91	2.01	2.09	2.15	2	3.33

- With the pair plots between parameters are:

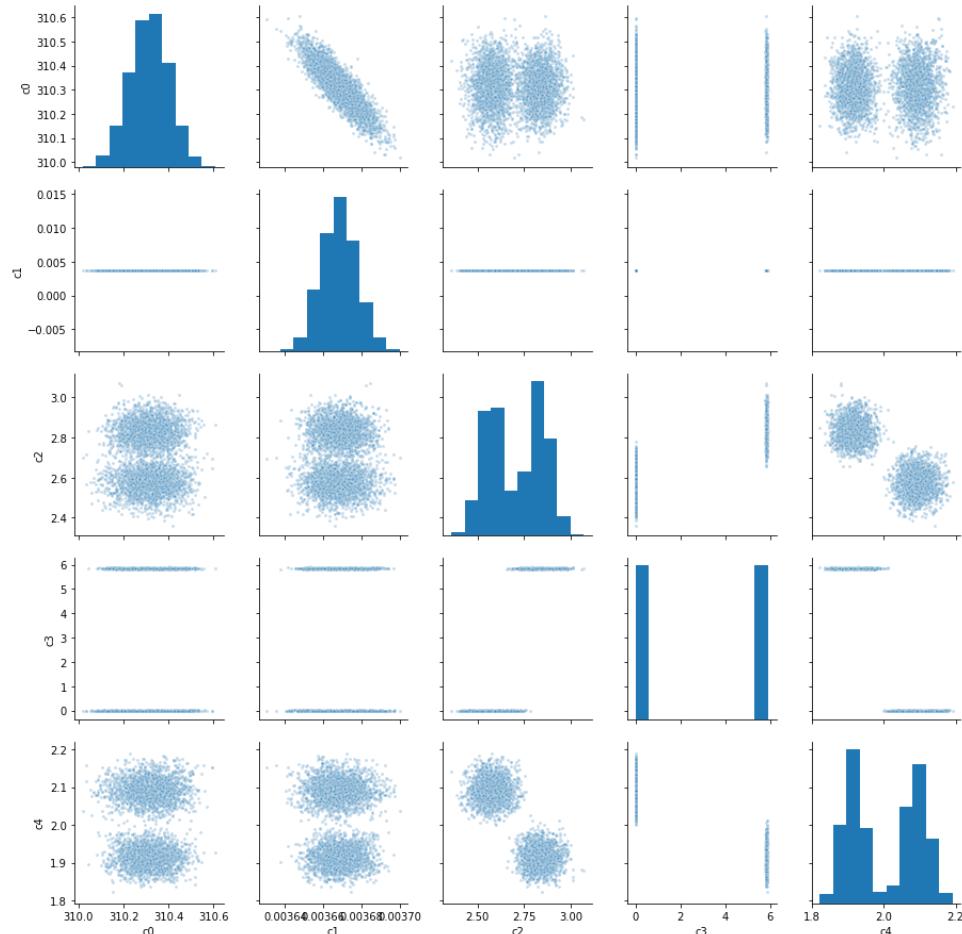


Fig 2: The pair plots visualize the relationship between the parameters plugging into the model.

- We can see that variable c_3 is causing the trouble as it is the clearest bimodal distributions. However, we realize that $\cos(x) = \cos(-x)$, which might be the reason for bimodal. As c_3 is

running freely and periodic in $[0, 2\pi]$ Hence, we narrow the region for c_3 to be in $[0, \pi]$.

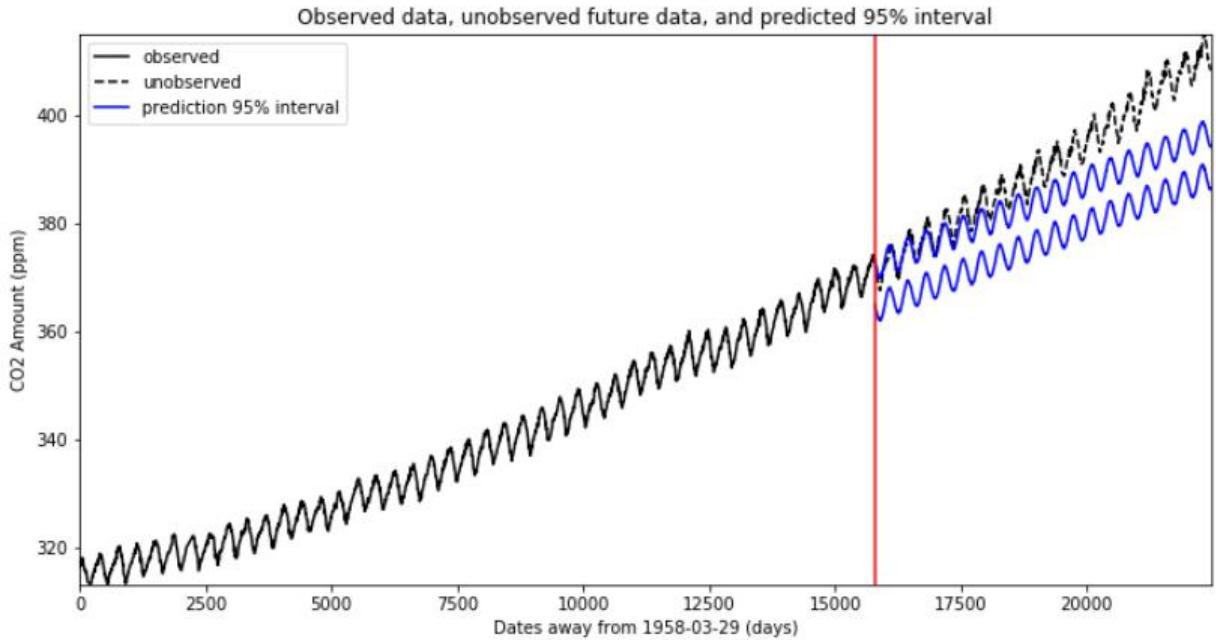


Fig 3: The observed data, unobserved data, and our predicted 95% interval.

Also, for the prediction, we can see that it is very far off, which might be because of the linearity. Hence, in the revised model, I will add a quadratic term.

Part 3: Model 2: Revision for the example model

$$p(x_t|\theta) = N\left(c_0 + c_1 t + c_5 t^2 + c_2 \cos\left(\frac{2\pi t}{365.25} + c_3\right), c_4^2\right)$$

- Changes:
 - o Narrow $c_3 \in [0, \pi]$
 - o Add quadratic term $c_5 t^2$
- Results: the model converges, with little correlated samples (based on the large n_eff values and Rhat values = 1).

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
c0	310.32	1.5e-3	0.09	310.13	310.25	310.32	310.38	310.5	3950	1.0
c1	3.7e-3	1.6e-7	1.0e-5	3.6e-3	3.7e-3	3.7e-3	3.7e-3	3.7e-3	3953	1.0
c2	2.57	1.3e-3	0.06	2.45	2.53	2.57	2.62	2.7	2263	1.0
c3	1.3e-3	2.4e-5	1.3e-3	2.3e-5	3.5e-4	8.5e-4	1.7e-3	4.9e-3	2935	1.0
c4	2.1	7.2e-4	0.03	2.03	2.07	2.09	2.12	2.16	2188	1.0

The prediction interval for unobserved data is:

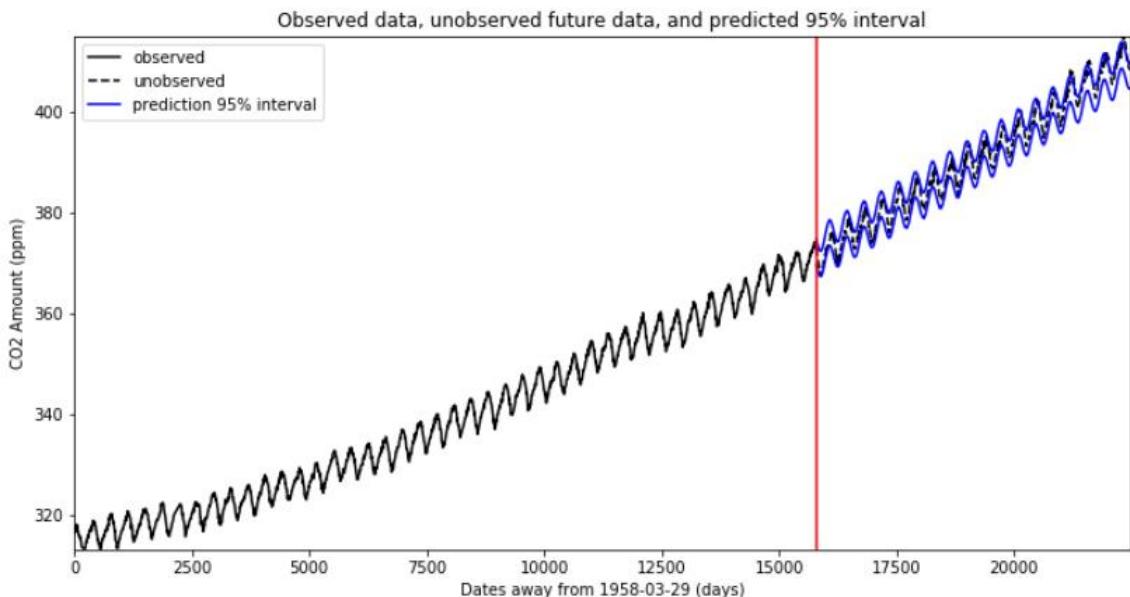


Fig 4: The observed data, unobserved data, and our predicted 95% interval

- We can see that this model fits the future unobserved data well; however, as the cosine function has a slightly curvy peak. When we look at the data, especially at the peaks, we observe that it is sharper than those of the predicted interval. Hence, I propose using the sawtooth wave function to recreate the sharp peaks.

Part 4: Model 3: Proposed model

$$p(x_t|\theta) = N\left(at^2 + bt + c + A \left(4 * \left| \left(\frac{t}{365.25} + \phi\right) \bmod 1 - 0.5 \right| - 1\right), \sigma^2\right)$$

- Changes:
 - o Using the sawtooth period function to achieve a sharp peak.
 - o Modify a to be distributed by a log-normal prior. The value seen in the model will be $\log(a)$, denoted by "log_a."
- Results: the model converges, with little correlated samples (based on the large n_{eff} values and Rhat values = 1).

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
A	3.41	1.0e-3	0.04	3.34	3.38	3.41	3.44	3.48	1294	1.0
phi	0.93	2.8e-5	1.7e-3	0.93	0.93	0.93	0.93	0.94	3668	1.0
log_a	-16.23	3.6e-4	0.01	-16.25	-16.23	-16.23	-16.22	-16.2	1186	1.0
b	2.2e-3	5.5e-7	1.8e-5	2.2e-3	2.2e-3	2.2e-3	2.2e-3	2.3e-3	1114	1.0
c	314.2	1.9e-3	0.07	314.07	314.16	314.2	314.25	314.33	1151	1.0
sigma	0.98	3.0e-4	0.01	0.95	0.97	0.98	0.99	1.0	2412	1.0

- The prediction for future data is:
(and the pair plots between the parameters)

We can see that at the peaks, we have the sharp peaks as we desired, and the intervals seem a bit smaller compared to model 2 (the revised example model), which indicates a slightly better model fit.

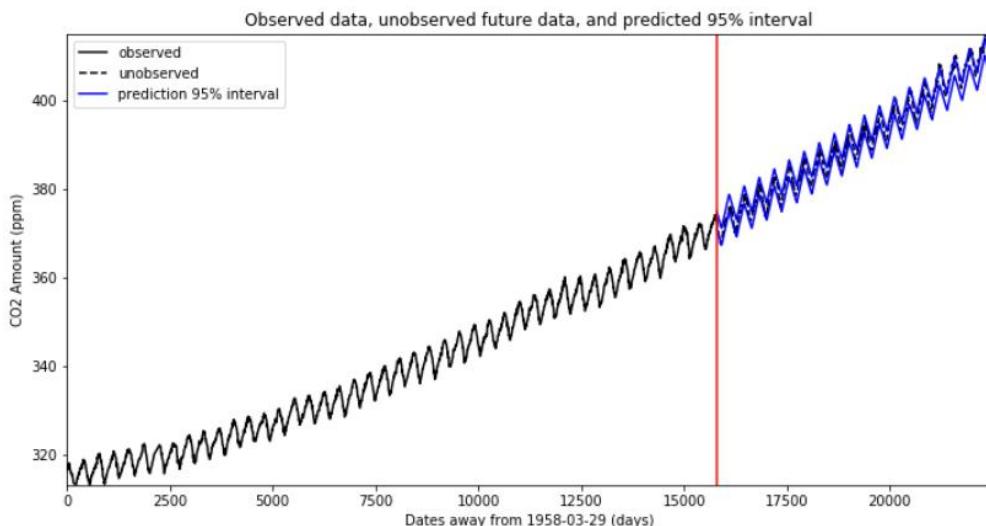


Fig 5: The observed data, unobserved data and our predicted 95% interval

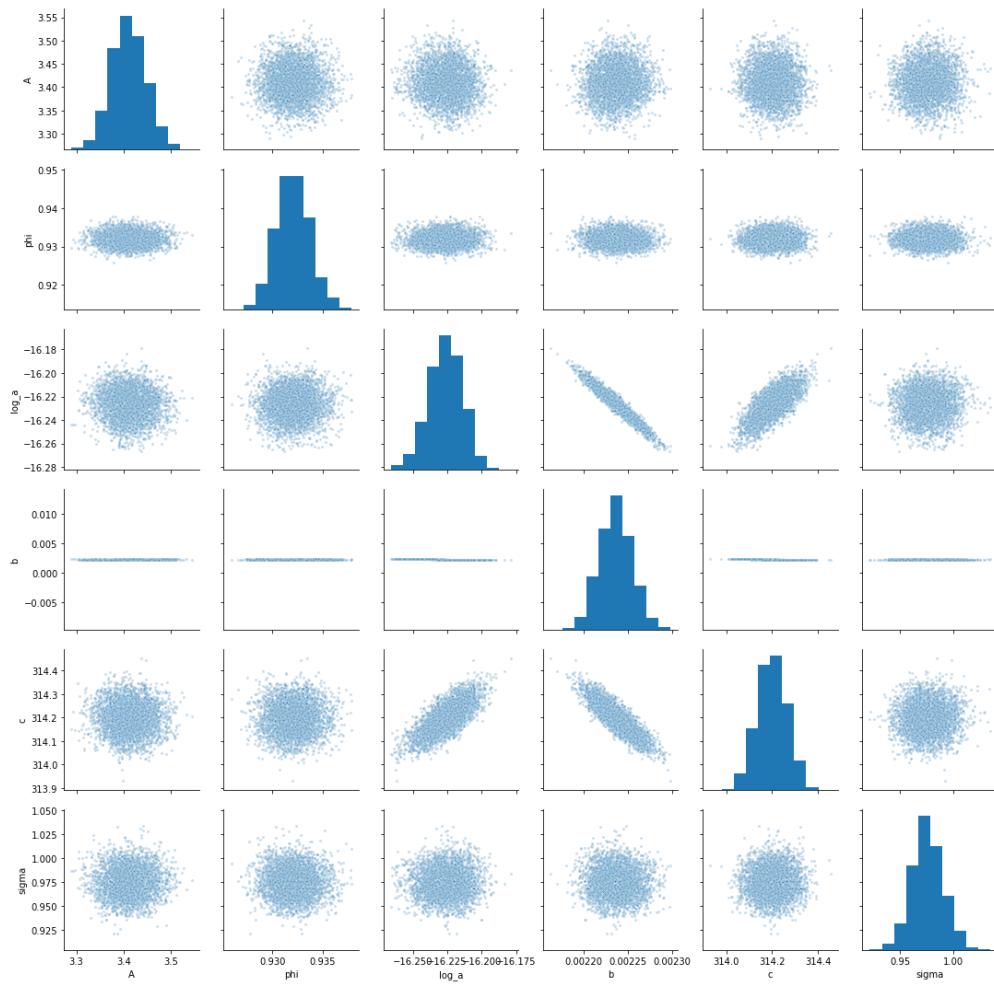


Fig 6: The pair plots visualize the relationship between the parameters plugging into the model.

Part 5: Model comparison:

Model 1 has many flaws: did not converge, and the linear component cannot describe the future data well.

Between model 2 (cosine periodic) and model 3 (sawtooth periodic):

Both models did an excellent job in predicting the future data: we see the prediction lines are very close to the unobserved data.

I prefer model 3 because:

- The shape of the peak of the data is sharp, which is more similar to the sawtooth than the cosine function.
- The prediction interval of model 3 is much narrower than of model 2: 100% of the intervals in model 3 is more restricted, which is a plus: indicates less uncertainty.
- Model 3 has lower predictive error: When we conduct mean squared errors on:
 - o The mean of the samples at time t vs. the real value at time t :
Model 3: error = 1.612156603044492
Model 2: error = 2.1612224291754356
 - o Every sample at the time t vs. the real value at time t :
Model 3: error = 2.5857090803350777
Model 2: error = 3.824635218251834

Model 3 has more accurate predictions: model 3 both have a lower predictive error when comparing the real data with the mean of the samples or with all the samples values themselves.

Furthermore, it has a much narrower interval (smaller $\Sigma |97.5\% \text{ interval} - 2.5\% \text{ interval}|$), which implies less uncertainty.

Hence, model 3 will be the model we use for the remaining of this project.

As we decided to move forward with model 3, we rerun the model with 100% of the data we have. Here is the graph demonstrate the predicted CO₂ level in the next 40 years.

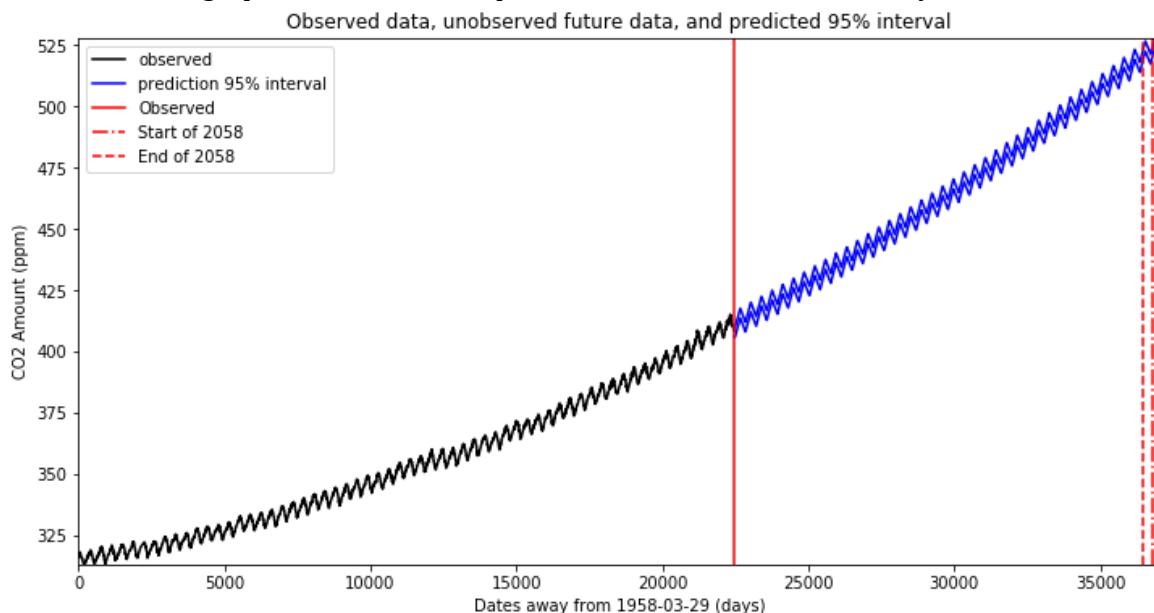


Fig 7: The observed data (current data) and our predicted 95% interval for future data.

Part 6: Predictions for the year 2058: the first measurement in 2058

Using our model, we have the predicted values for CO₂ Amount in the first measurement in 2058:

- With the mean at 519.5426720521242 (ppm)
- The median at 519.5288075969152 (ppm)
- Standard deviation at 1.073557430664101
- The 95% interval at: [517.4788034343148, 521.6107600881326]

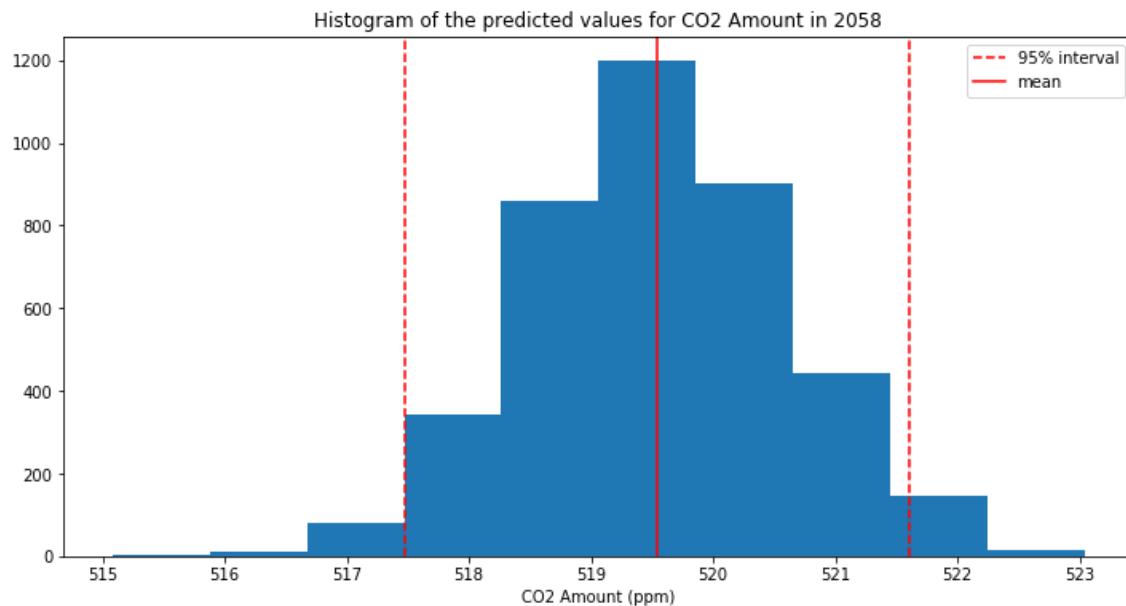


Fig 8: The histogram of the estimates of the CO₂ Amount in the 1st measurement in 2058

The graph here is drawn from 4000 samples data. There is a possibility that the samples are underestimating the values (which discussed later in part 8 and 9).

Part 7: The probability of CO₂ level reaching 450 ppm

Using the samples from our predicted values, we evaluate the probability of reaching 450 ppm by calculating among 4000 samples for a specific date, what proportion is more significant than 450 ppm.

These probabilities are dependent on the samples, which might not be fully reflective of future scenarios. Probability = 1 does not necessarily 100% happening; it only indicates that 100% of the samples satisfy the conditions.

We consider two thresholds: 0.5 and 1.0

The dates are only estimating and did not account for uncertainty. Uncertainty will be discussed in Part 8 &9 (shortcomings and recommendations).

- The first predicted date with the probability of reaching 450 ppm larger than 0.5 is "03-18-2034".
- The first date with the probability of reaching 450 ppm equal to 1.0 is "04-07-2035".
- The first date with the probability of everyday onward reaching 450 ppm larger than 0.5 is "11-17-2035".
- The first date with the probability of everyday onward reaching 450 ppm equal 1.0 is "11-29-2036".

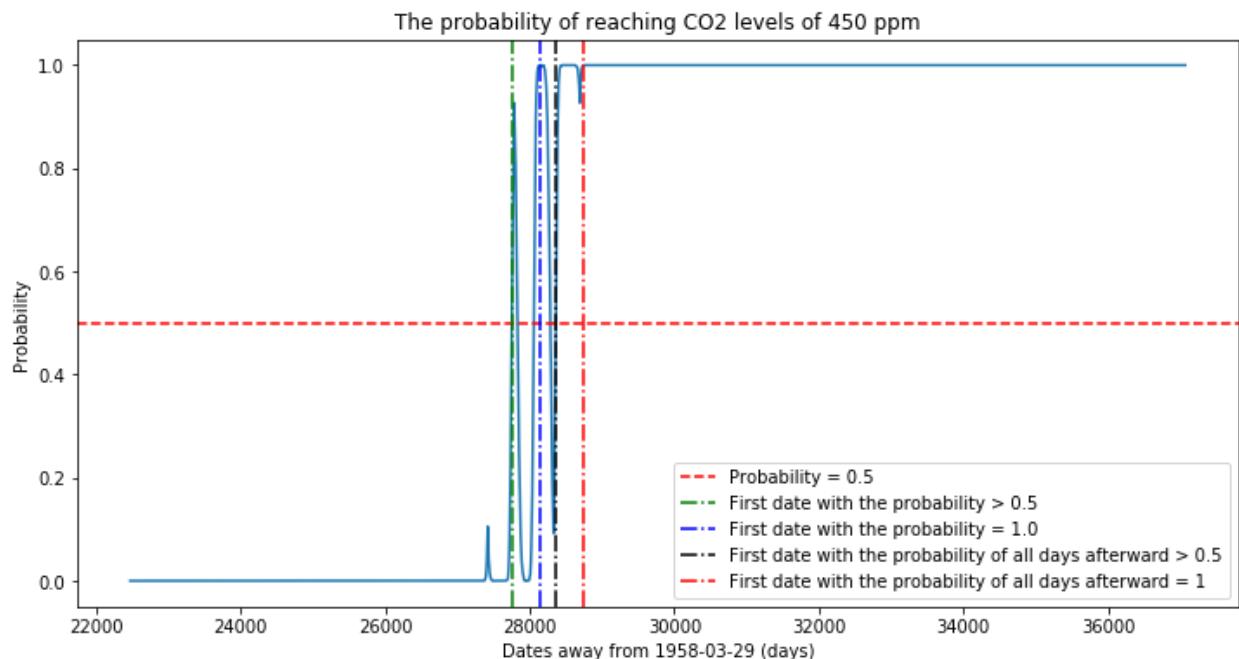


Fig 9: The probability that the estimates reaching levels of 450 ppm.

Part 8: Shortcomings:

Shortcomings

- a. The process of designing the model
 - i. Assumptions
 - ii. Priors
 - iii. Parameters
 - iv. Observed vs. Unobserved variables
- b. Test statistics

Part 8.a.: The process of designing the model

Part 8.a.i:

Assumption 1: Assuming the data has some periodic component by year (365.25). This assumption is reasonable because the year (contains the seasons), which directly affects the human. Also, this assumption helps the model because if the data is periodic by 1 year $\lambda = \frac{1}{365.25}$, it can be periodic by k years, with k is an integer ($\lambda = \frac{1}{365.25k}$), which might create divergences in stan sampling.

Assumption 2: Our model assumes that the increasing factor is quadratic. The reason is that we try to fit a polynomial factor as the increasing factor. The linear fails with model 1, which inspires us to use quadratic, which resulted in a very close fit model. We did not model higher-order polynomials or exponential because we suspect the growth rate is not that fast.

Assumption 3: We assume the noise factor is a Gaussian mean with a fixed standard deviation σ .

Assumption 4: All the priors' assumptions.

Part 8.a.ii. + 8.a.iii. + 8.a.iv.:

$A \sim N(0,5)$: the range of the periodic factor will be from $[-A, A]$. We choose $N(0,5)$ because the CO₂ measurements periodic fluctuation should not be too big. However, as we are unsure, we still keep a pretty big standard deviation.

For ϕ (the initial state), because we take modular 1, we need to rescale ϕ to $[0,1]$ by $\frac{\arctan\frac{\phi_1}{\phi}}{2\pi} + 0.5$: $\arctan\frac{\phi_1}{\phi}$ will take values to $[-\pi, \pi]$, then the division and addition will transform this to $[0,1]$.

For the quadratic component: b will be the main components for the increase, a will create the curve and increase when t gets larger, and c will represent the initial value at the beginning of the recording period. $at^2 + bt + c$

For a (quadratic term): we want a to be small (as the time component is huge), and the change of the CO₂ should not be too significant and center around a millionth. Hence, we model $\log(a) \sim N(-16, 6)$.

For b (linear term): we want b to be the primary increasing factor: $b \sim N(0, 0.1)$.

For c (constant term): c is the initial state when we start recording the data. We use a broad term to demonstrate uncertainty and allow more freedom for the model. $c \sim N(300, 100)$

σ (the noise): we are unsure, so we choose a broad range: $\sigma \sim \text{Gamma}(1, 0.3)$

As we believe the amount is increasing, we set the limit for all parameters to have lower = 0:
 $A, \sigma, a, b, c \geq 0$

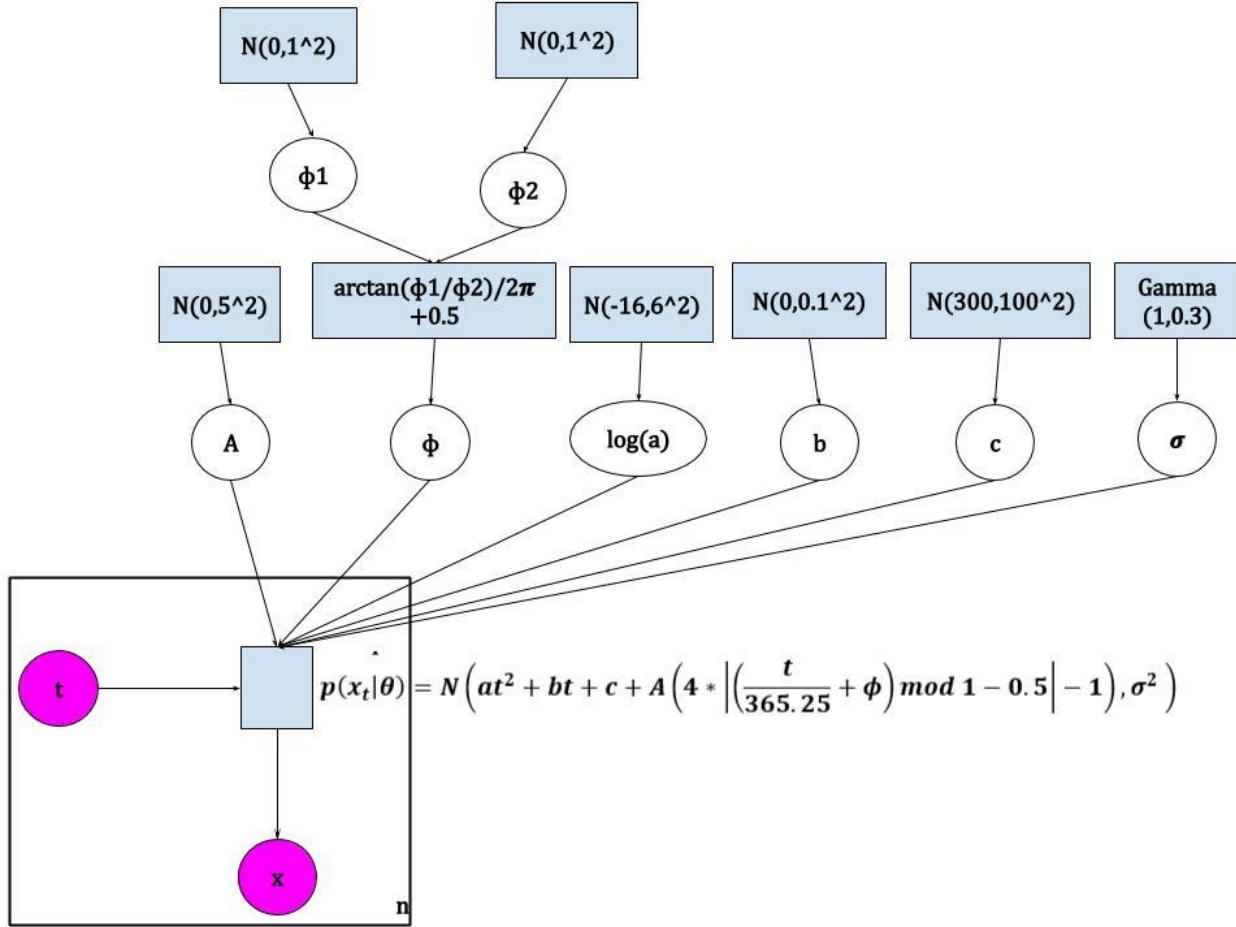


Fig 10: The factor graph of the sawtooth-quadratic model (model 3)

Part 8.b. test Statistics: As discussed above, we use the first 70% of the data (2197 values) as data for computing the posterior distributions over the parameters above using Stan. We use the other 30% of the data (the next 942 values) to test our prediction of future function values.

Our test statistic is the mean of future data: 942 last data points in the dataset.

The reason is that: we suspect the model, even though overall fit the future data very well (based on the plot), fails to predict the few anomalies in the original data fully. For example, when we zoom in the first 600 data points in the test future data, we can see at those peaks, our real data is slightly off compared to the predicted interval.

As the mean is sensitive to the outliers, we would expect the mean for the real future data to be in the upper or bottom percentile of the test statistics results of the sampled data.

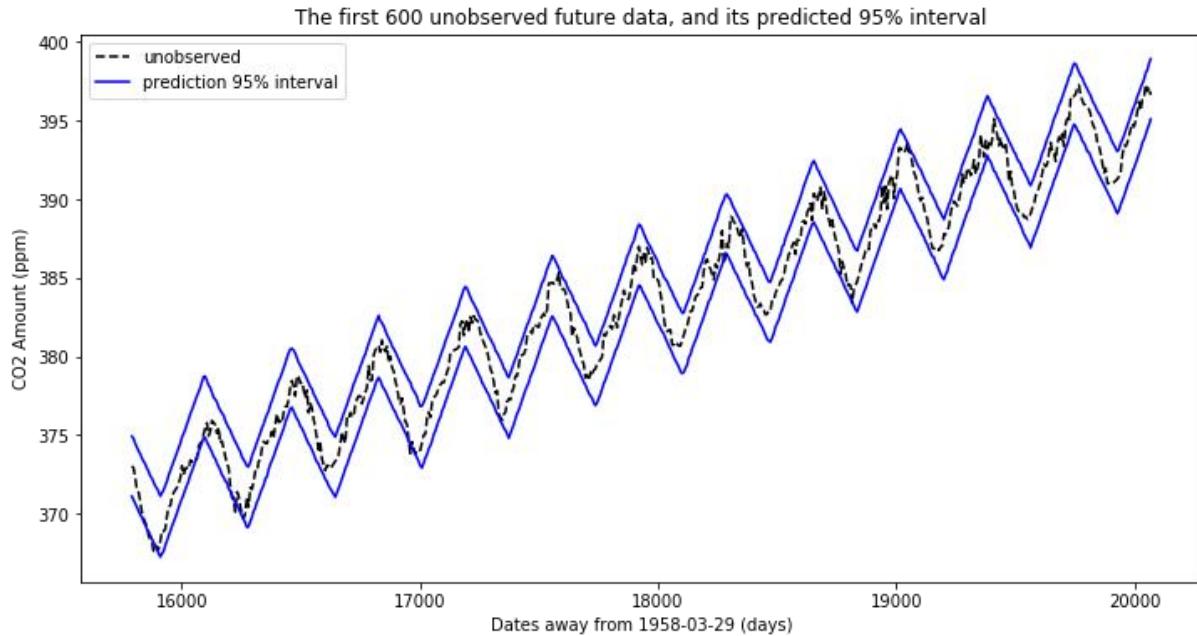


Fig 11: The zoom-in version for the first 600 datapoints to investigate the performance of the model

We generate 4000 sampled datasets (each dataset is an estimate for the 30% last unobserved data)

We then computed the mean of those 4000 datasets and compared to the mean of the real unobserved data.

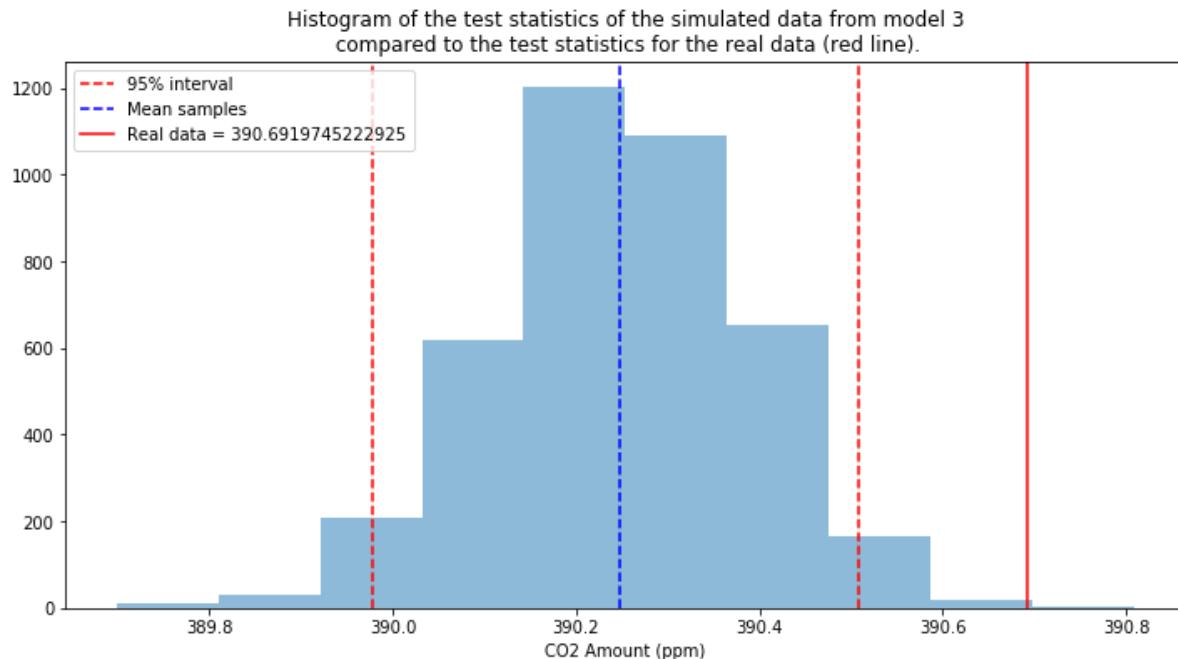


Fig 12: The histogram of the test statistics simulated from model 3.

The percentage that samples result larger than the real statistic value: 0.1 %

Hence, the p-value: 0.999: which is more significant than 0.95, meaning that the model failed to sufficiently predict the unobserved data with respect to this test statistics.

Looking at the histogram, we can infer that the predictive results by our model tend to underestimate the real data. Model 2 faced the same failure in this test statistic (p-value = 0.9928).

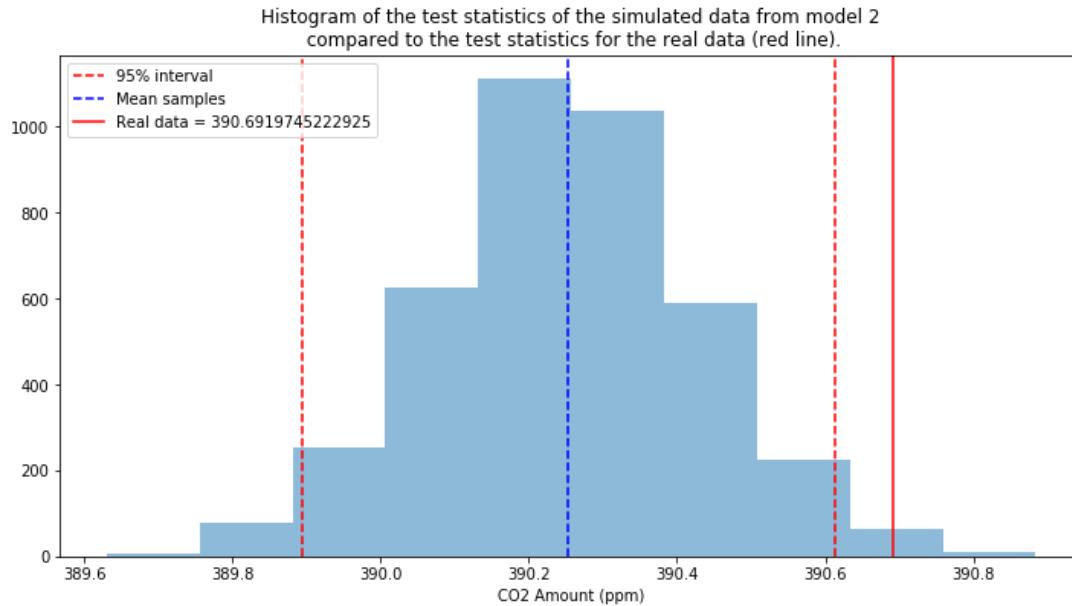


Fig 12: The histogram of the test statistics simulated from model 3.

To investigate further what might be the reason, we measure the errors overtime:
(we did not plot the 95% interval because it would look too clustery)

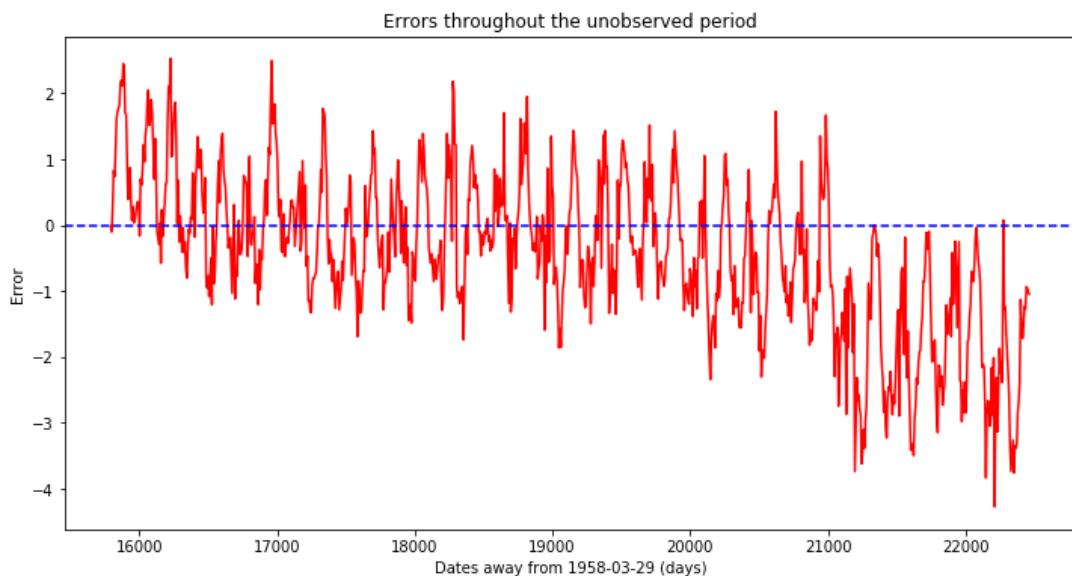


Fig 13: The average errors throughout the unobserved period.

As we can see, the error starts to underestimate the real values as time goes further. This could mean that our quadratic component did not increase fast enough to fit the data. We could try an exponential or a higher-order polynomial fit. However, for the scope of this project, and with the limited data, I would not add more parameters, which might overfit the data.

This graph would also raise concerns regarding our estimates for 2058 data and high-risk CO₂ levels, as those both are far in the future. Hence, we might underestimate even more for those values (2058 could have a higher CO₂ level, while the date with high probability for high-risk CO₂ could come faster). This will be further discussed in Part 9.

The set of the errors between all samples and the real data is very broad: with the median at -0.4 and the 95% interval at [-3.63686949, 2.46959212], and especially getting more negative while t gets larger.

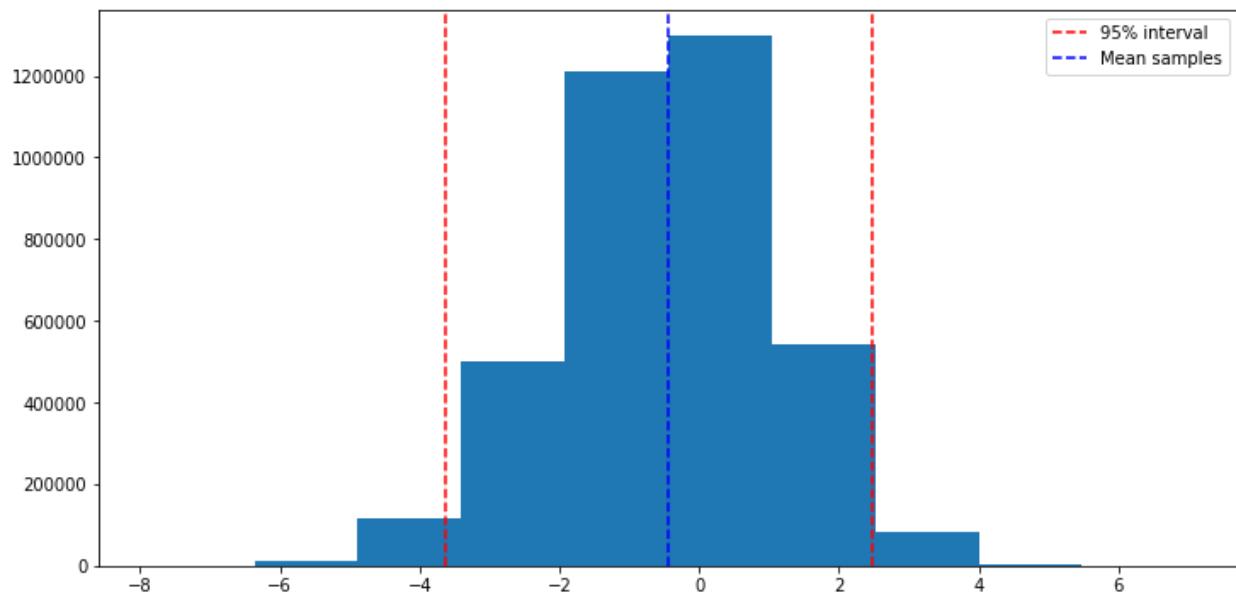


Fig 14: The histogram of all errors throughout the unobserved period.

Part 9: Discussion and Recommendations

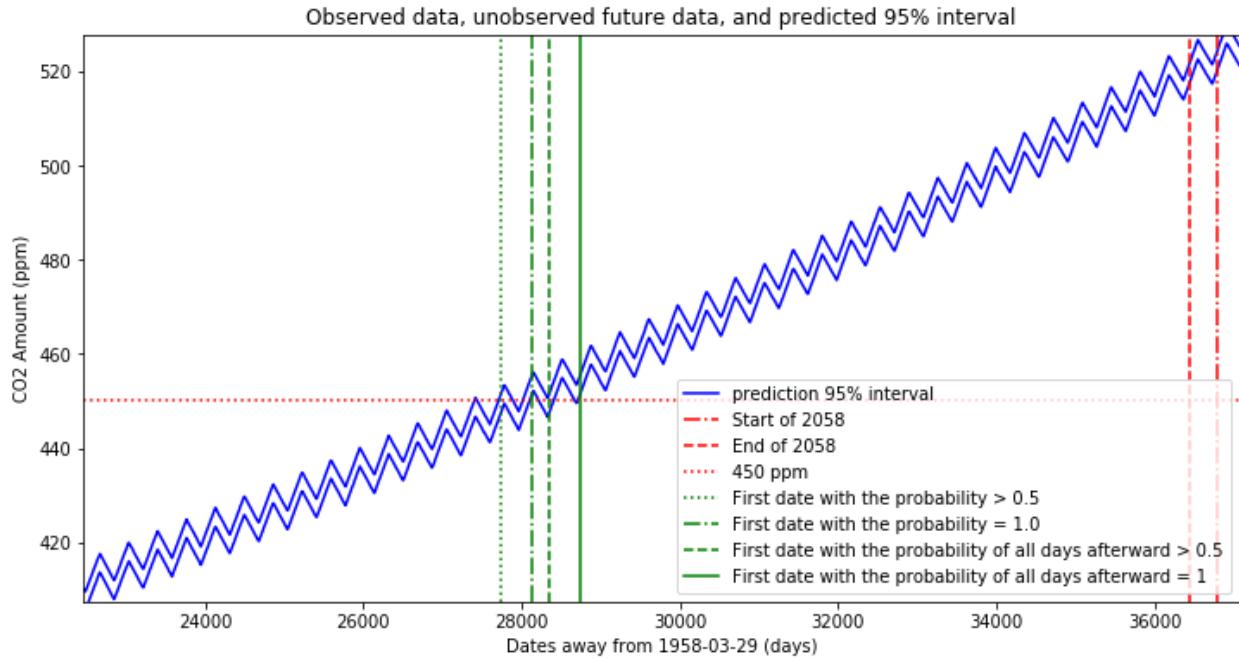


Fig 15: The plot of important components of this project

For the CO₂ modeling, we create a model with an increasing quadratic trend, a sawtooth component to represent seasonal factors and a Gaussian noise center at 0 with standard deviation σ .

The model fit the data quite well, with an average 1.6 ppm error when we subset the data into 70% train and 30% to validate. Overall, the model has a very closed estimation of the real data.

Notable drawback: the model tends to underestimate the real data as time gets further from the initial stage. Hence, the prediction for the year 2058 could be an underestimate (the values could be higher than what the model suggested). Also, the forecast for the high-risk CO₂ level could be earlier than what the model indicated.

Recommendations:

- For the high-risk CO₂ level (part 7), each probability is determined by 4000 samples from Stan sampling. If we want to have uncertainty for the probability estimate, I suggest running the Stan sampling multiple times to include uncertainties and a more accurate date range.
- Account for the possible underestimates while doing estimates: the current error is around -0.4 (the median of the errors set), with a confidence interval [-3.6, 2.5], with the errors getting more negative when t get larger.
- With the error sets data (972 time-values, each has 4000 error data points), try to figure a model to explain the systematic error (as it seems to be a systematic decrease).
- Or try fitting the model with either exponential or higher-order polynomials for the trend component to check if a more accurate model could be found.
- Save the environment!

HCs Applications

#communicationdesign: For this final project, I focus on demonstrating events and information through multiple data visualizations, as graphs are more intuitive to understand and interpret, both for statisticians and the general audience. This project is a good use of #communicationdesign as I take into consideration the audience and the representations of the deliverable. Also, I adhere to the cognitive principles of the communications design: reduce cognitive workload by using figures, each figure focuses on delivering one result, and all figures are clearly labeled and explained.

#testability: For all models, I focus on validating the advantage and disadvantages. I subset the data into train and validate data and conduct appropriate tests to analyze and compare with other models on different metrics (predictive error, the width of confidence intervals). Also, I design a test statistic that addresses the possible shortcomings of the models (the mean), which addresses the underestimation when t becomes large. This is a good use of #testability because I design testable models and test them thoroughly.

#casestudy: Using the data for this specific case: CO₂ level recorded at Mauna Loa Observatory, I try to model and extract as labeled insights as possible on the behaviors of the CO₂ level. I identify the suitable components to explain the data, describe the limitations, and provide recommendations to extend the study. This is a good use of #casestudy as I utilize a set of data on many aspects and infer trends when applicable

#selfawareness: While doing this project, I acknowledge that I might be subject to Dunning-Kruger effect (unable to recognize my incompetence). Hence, I try to start the project early and go through multiple stages of questions and feedback with Professor. Scheffler, which helps mitigate my weaknesses. This strategy works exceptionally well with me as I tend to procrastinate, which, through trying to discuss in Office Hour, helps me to motivate myself to work harder to have a more insightful conversation.

References:

Sampling Station Records - Scripps CO₂ Program - Mauna Loa Observatory, Hawaii. (n.d.).
Retrieved December 2019, from https://scrippsc02.ucsd.edu/data/atmospheric_co2/mlo.html.

```
1 # import packages
2
3 import pystan
4 import numpy as np
5 from scipy import stats
6 import matplotlib.pyplot as plt
7 import pandas as pd
8 import seaborn
9 import datetime
10 from datetime import timedelta
11 %matplotlib inline
12 np.random.seed(2019)
```

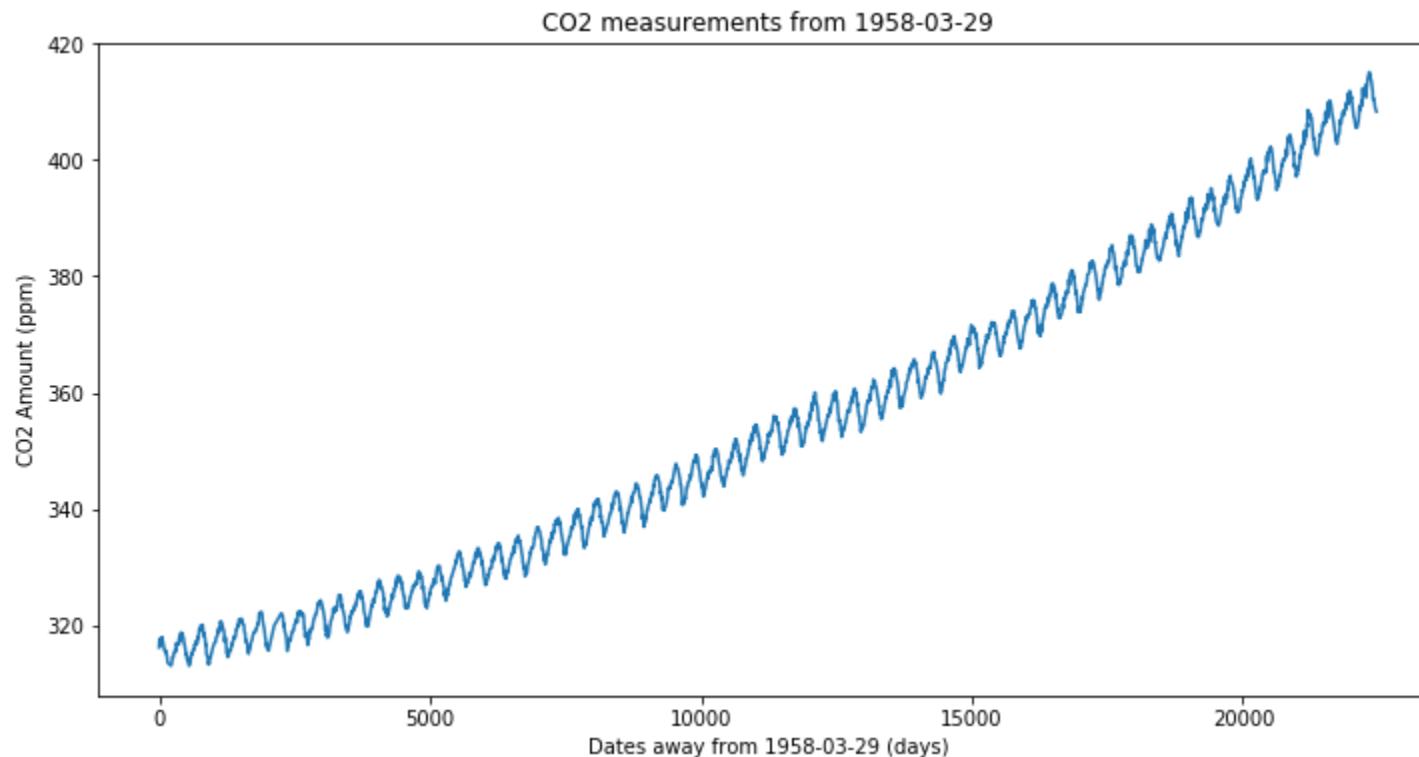
▼ Data Processing

```
1 data = pd.read_csv("https://docs.google.com/spreadsheets/d/e/2PACX-1vTX5nZNL-okv1cEdQ00ue57BEI6ARhDZwMJWiUyD4xupTcIx0aEI75UTxp8P9zx")
1 data.columns = ["Date", "Amount"]
1 max(data["Date"])
↳ '2019-09-28'
1 # subtract the date values from the first date
2 time = []
3 for i in range(len(data)):
4     time.append((datetime.datetime.strptime(data["Date"][i], "%Y-%m-%d") - datetime.datetime.strptime(data["Date"][0], "%Y-%m-%d")).days)
1 data["Time"] = time
1 data.head()
↳
```

	Date	Amount	Time
0	1958-03-29	316.19	0
1	1958-04-05	317.31	7
2	1958-04-12	317.69	14
3	1958-04-19	317.58	21
4	1958-04-26	316.48	28

```
1 plt.figure(figsize = (12,6))
2 plt.title('CO2 measurements from 1958-03-29')
3 plt.plot(data["Time"], data["Amount"], label = "Real data")
4 plt.xlabel("Dates away from 1958-03-29 (days)")
5 plt.ylabel("CO2 Amount (ppm)")
6 plt.show()
```

➡



▼ Simple model: Linear + cos periodic

```
1 # The Stan model. Running this cell compiles the Stan model, which takes
2 # some time.
3 # Stan gets to see the first n data values only
4
5 stan_code = """
6 data {
7     int<lower=0> n;           // The number of data
8     int<lower=0> n_future;    // the number of predictions
9     real x[n];               // The data
10    real t[n];                // time data
11    real t_future[n_future]; // future time data
12 }
```

```

14 parameters {
15     real<lower = 0> c0;
16     real<lower = 0> c1;
17     real<lower = 0> c2;
18     real<lower = 0> c3;
19     real<lower = 0> c4;
20 }
21
22 model {
23     for(i in 1:n) {
24         x[i] ~ normal(
25             c0 + c1*t[i] + c2*cos(2*pi()*t[i]/365.25 + c3),
26             c4);
27     }
28 }
29
30 // Generate the predicted function values for the next n_future steps.
31 generated quantities {
32     real x_future[n_future];
33     for(i in 1:n_future) {
34         x_future[i] = normal_rng(
35             c0 + c1*t_future[i] + c2*cos(2*pi()*t_future[i]/365.25 + c3),
36             c4);
37     }
38 }
39 """
40
41 stan_model = pystan.StanModel(model_code=stan_code)

```

↳ INFO:pystan:COMPILE THE C++ CODE FOR MODEL anon_model_ebe2c6f5518b031d8f01ff4ffba65ddc NOW.

```

1 # subsetting the data
2 n = len(data)*7//10
3 n_future = len(data) - n
4
5 # stan data
6 stan_data = {
7     'n': n,

```

```

8     'n_future': n_future,
9     'x': data["Amount"][:n],
10    't': data["Time"][:n],
11    "t_future" : data["Time"][n:]}
12
13 # Run Hamiltonian Monte Carlo using Stan with 4 Markov chains, a 1000 step
14 # warm-up phase and 1000 step sampling phase for each chain. The warm-up samples
15 # are discarded, so we are left with 4 x 1000 = 4000 samples.
16
17 # This cell will take a minute or two to run.
18
19 parameters = ['c0', 'c1', 'c2', 'c3', 'c4']
20
21 results = stan_model.sampling(data=stan_data)
22 print(results.stansummary(pars=parameters))
23 samples = results.extract()

```

↳ WARNING:pystan:n_eff / iter below 0.001 indicates that the effective sample size has likely been overestimated
 WARNING:pystan:Rhat above 1.1 or below 0.9 indicates that the chains very likely have not mixed
 Inference for Stan model: anon_model_ebe2c6f5518b031d8f01ff4ffba65ddc.
 4 chains, each with iter=2000; warmup=1000; thin=1;
 post-warmup draws per chain=1000, total post-warmup draws=4000.

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
c0	310.31	1.4e-3	0.09	310.15	310.25	310.31	310.37	310.48	3822	1.0
c1	3.7e-3	1.6e-7	9.3e-6	3.7e-3	3.7e-3	3.7e-3	3.7e-3	3.7e-3	3546	1.0
c2	2.7	0.09	0.14	2.47	2.58	2.7	2.83	2.94	2	2.39
c3	142.72	148.08	209.46	4.9e-5	7.9e-4	34.31	307.12	502.24	2	1.6e4
c4	2.0	0.07	0.1	1.87	1.91	2.01	2.09	2.15	2	3.39

Samples were drawn using NUTS at Wed Dec 18 13:51:30 2019.
 For each parameter, n_eff is a crude measure of effective sample size,
 and Rhat is the potential scale reduction factor on split chains (at
 convergence, Rhat=1).

```

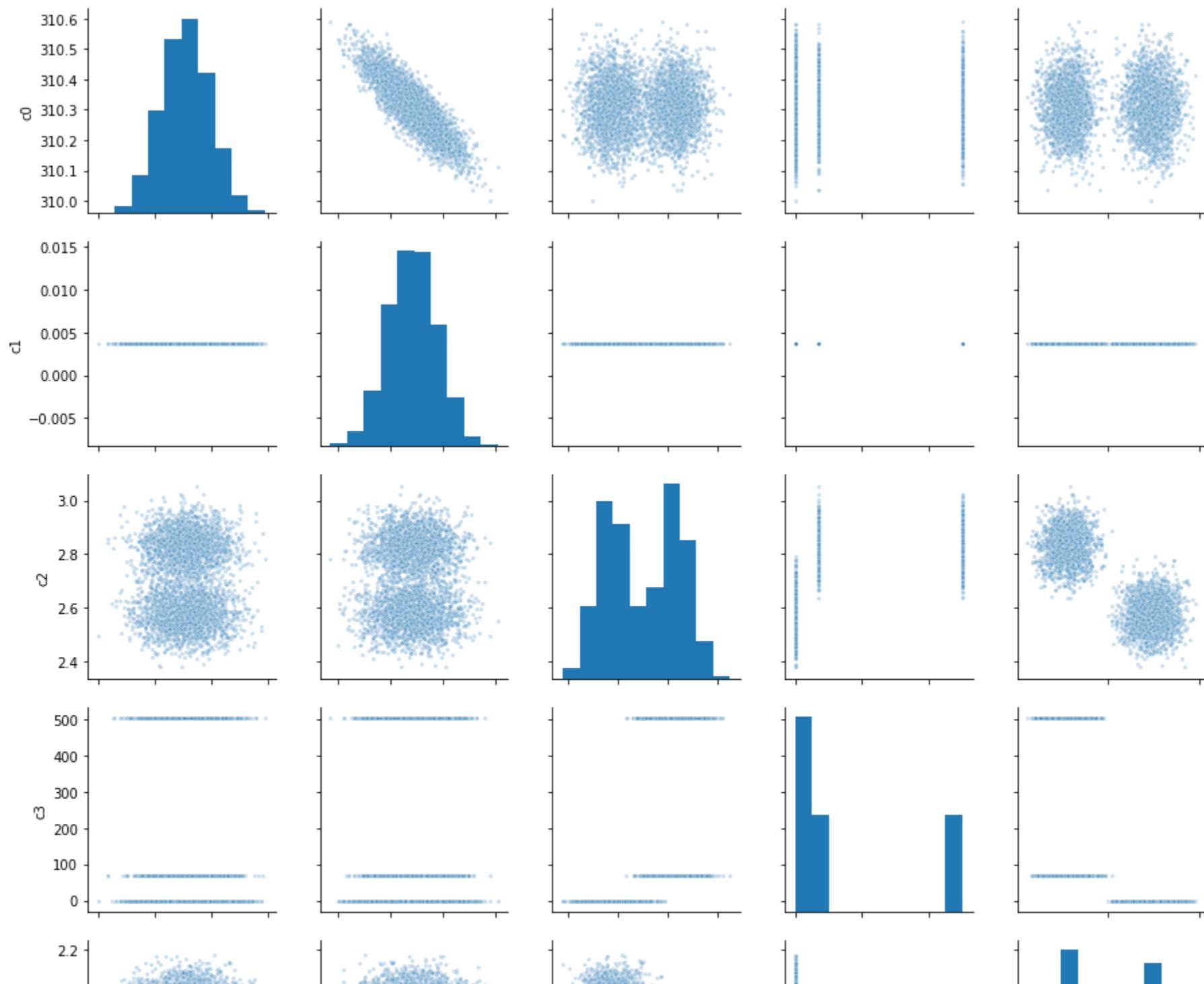
1 # plot the pair plot
2
3 df = pd.DataFrame(
4     data=np.transpose([samples[param] for param in parameters]),
5     columns=parameters)

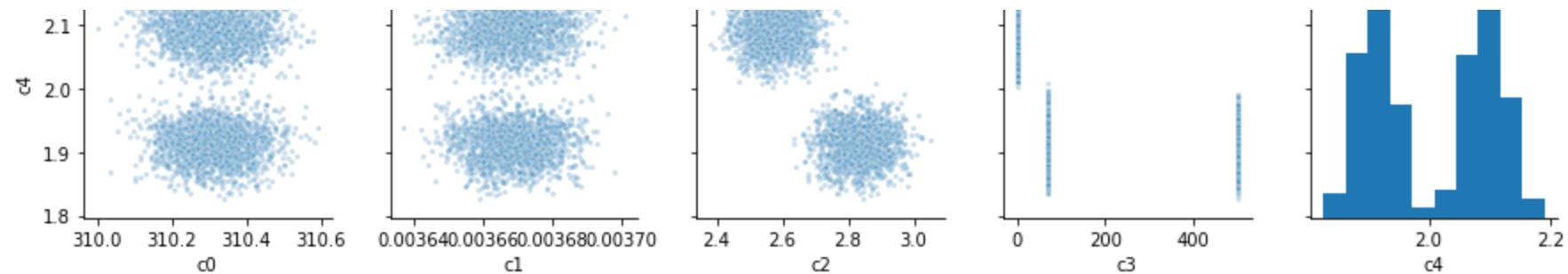
```

<https://colab.research.google.com/drive/1KWIUWOi9tH63n6MPaxxKX8InGOQCMSk?authuser=1#scrollTo=KGO8Kw2-v8DS&printMode=true>

```
6 seaborn.pairplot(df, height=2.5, plot_kws={'marker': '.', 'alpha': 0.25})  
7 plt.show()
```





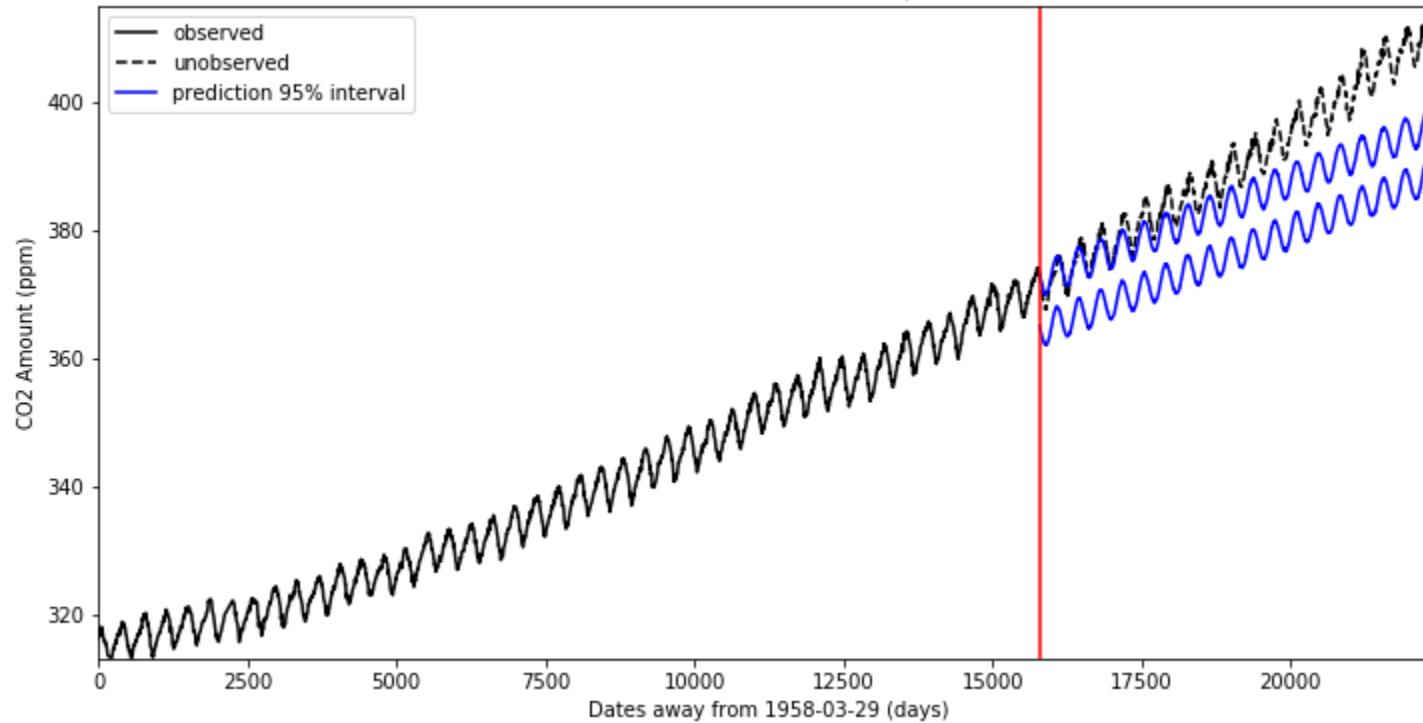


```
1 # draw predictions
2 prediction = samples['x_future']
3
4 # Compute 95% interval of the predicted values
5 prediction_interval = np.percentile(prediction, axis=0, q=[2.5, 97.5])
6
7 # Plot mean and 95% interval of predictions
8 plt.figure(figsize=(12, 6))
9 # observed data
10 plt.plot(data["Time"][:n], data["Amount"][:n], 'k-', label='observed')
11 #unobserved data
12 plt.plot(data["Time"][n:n+n_future], data["Amount"][n:n+n_future], 'k--', label='unobserved')
13
14 # set the grid
15 plt.xlim(0, data["Time"][n+n_future-1])
16 plt.ylim(
17     min(min(prediction[0,:]), min(data["Amount"])),
18     max(max(prediction[1,:]), max(data["Amount"])))
19
20 # confidence interval
21 plt.plot(
22     data["Time"][n:],
23     prediction_interval[0,:],
24     'b-', label='prediction 95% interval')
25 plt.plot(
26     data["Time"][n:],
27     prediction_interval[1,:],
28     'b-')
```

```
29  
30 plt.axvline(data["Time"][n - 1], color='red')  
31 plt.xlabel("Dates away from 1958-03-29 (days)")  
32 plt.ylabel("CO2 Amount (ppm)")  
33 plt.legend()  
34 plt.title('Observed data, unobserved future data, and predicted 95% interval')  
35 plt.show()
```



Observed data, unobserved future data, and predicted 95% interval



1

▼ Revised model: quadratic + cos periodic

```
1 # The Stan model. Running this cell compiles the Stan model, which takes  
2 # some time.  
3 # Stan gets to see the first n data values only  
4
```

```
5 stan_code = """
6 data {
7     int<lower=0> n;           // The number of data
8     int<lower=0> n_future;   // the number of predictions
9     real x[n];              // The data
10    real t[n];              // time data
11    real t_future[n_future]; // future time data
12 }
13
14 parameters {
15     real<lower = 0> c0;
16     real<lower = 0> c1;
17     real<lower = 0> c2;
18     real<lower = 0, upper = pi()> c3;
19     real<lower = 0> c4;
20     real<lower = 0> c5;
21 }
22
23 model {
24     for(i in 1:n) {
25         x[i] ~ normal(
26             c0 + c1*t[i] + c5*t[i]*t[i] + c2*cos(2*pi()*t[i]/365.25 + c3),
27             c4);
28     }
29 }
30
31 // Generate the predicted function values for the next n_future steps.
32 generated quantities {
33     real x_future[n_future];
34     for(i in 1:n_future) {
35         x_future[i] = normal_rng(
36             c0 + c1*t_future[i] + c5*t_future[i]*t_future[i] + c2*cos(2*pi()*t_future[i]/365.25 + c3),
37             c4);
38     }
39 }
40 """
41
42 stan_model = pystan.StanModel(model_code=stan_code)
```

↳ INFO:pystan:COMPIILING THE C++ CODE FOR MODEL anon_model_84f7f23059da5ade2b15b1965e6fdfad NOW.

```
1 # subsetting the data
2 n = len(data)*7//10
3 n_future = len(data) - n
4
5 stan_data = {
6     'n': n,
7     'n_future': n_future,
8     'x': data["Amount"][:n],
9     't': data["Time"][:n],
10    "t_future" : data["Time"][n:]}
11
12 # Run Hamiltonian Monte Carlo using Stan with 4 Markov chains, a 1000 step
13 # warm-up phase and 1000 step sampling phase for each chain. The warm-up samples
14 # are discarded, so we are left with 4 x 1000 = 4000 samples.
15
16 # This cell will take a minute or two to run.
17
18 parameters = ['c0', 'c1', 'c2', 'c3', 'c4', 'c5']
19
20 results = stan_model.sampling(data=stan_data)
21 print(results.stansummary(pars=parameters))
22 samples = results.extract()
```

↳

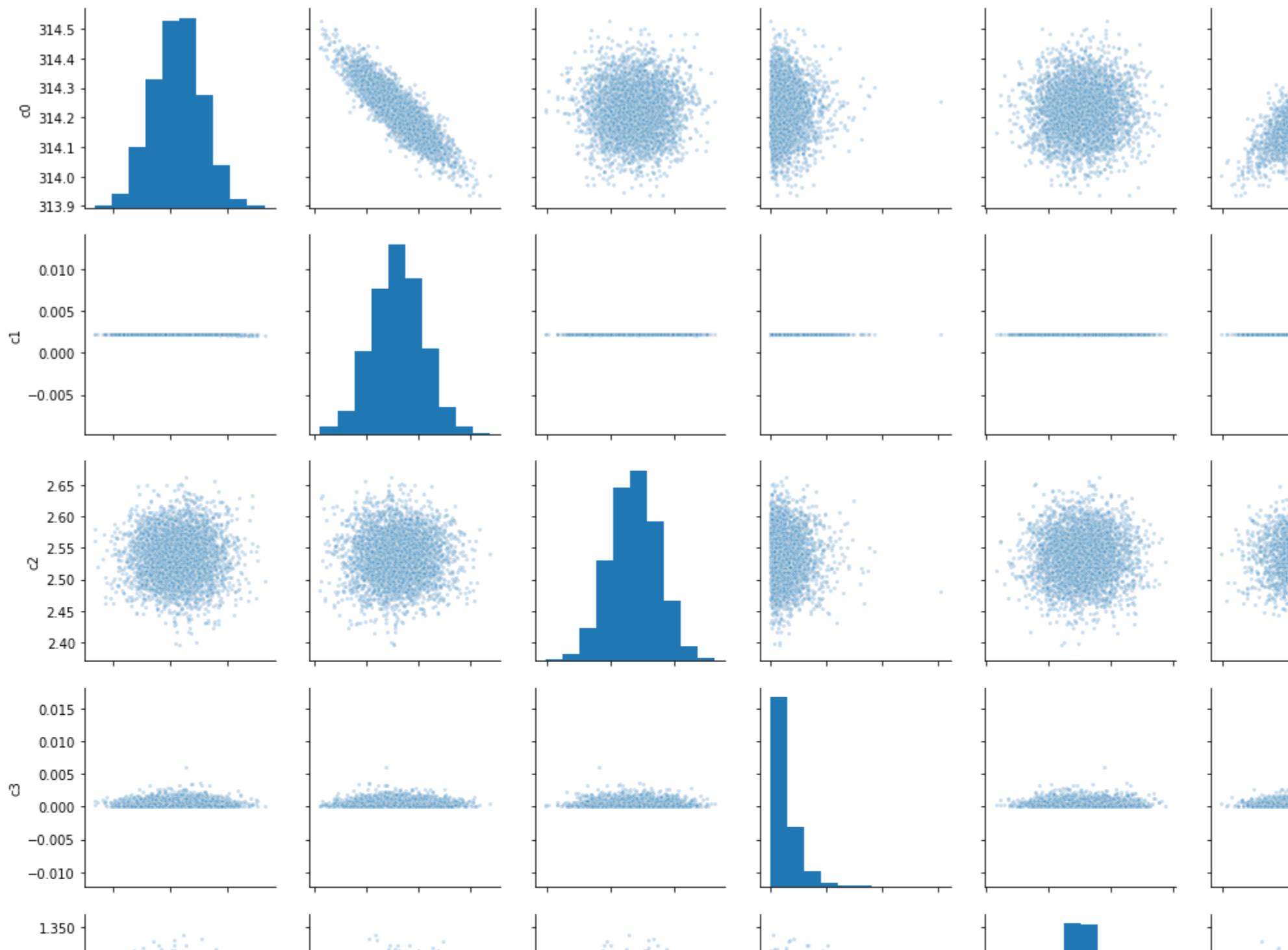
Inference for Stan model: anon_model_84f7f23059da5ade2b15b1965e6fdfad.
 4 chains, each with iter=2000; warmup=1000; thin=1;
 post-warmup draws per chain=1000, total post-warmup draws=4000.

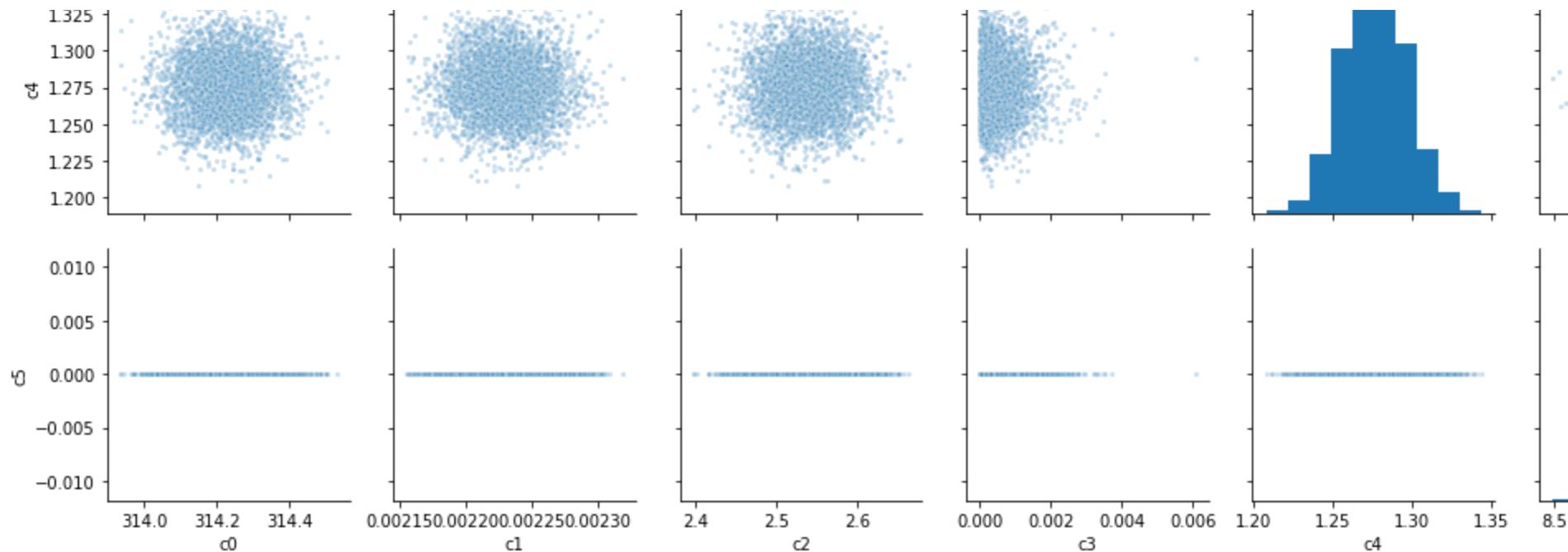
	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
c0	314.23	1.8e-3	0.09	314.06	314.17	314.23	314.28	314.39	2307	1.0
c1	2.2e-3	5.8e-7	2.4e-5	2.2e-3	2.2e-3	2.2e-3	2.2e-3	2.3e-3	1763	1.0
c2	2.54	6.7e-4	0.04	2.46	2.51	2.54	2.56	2.61	3226	1.0
c3	5.0e-4	8.4e-6	4.7e-4	1.3e-5	1.6e-4	3.6e-4	6.8e-4	1.7e-3	3097	1.0
c4	1.28	3.7e-4	0.02	1.24	1.26	1.28	1.29	1.32	2875	1.0
c5	9.0e-8	3.5e-11	1.5e-9	8.7e-8	8.9e-8	9.0e-8	9.1e-8	9.3e-8	1715	1.0

Samples were drawn using NUTS at Wed Dec 18 13:59:27 2019.
 For each parameter, n_eff is a crude measure of effective sample size,
 and Rhat is the potential scale reduction factor on split chains (at
 convergence, Rhat=1).

```
1 # pair plots
2
3 df = pd.DataFrame(
4     data=np.transpose([samples[param] for param in parameters]),
5     columns=parameters)
6 seaborn.pairplot(df, height=2.5, plot_kws={'marker': '.', 'alpha': 0.25})
7 plt.show()
```







```

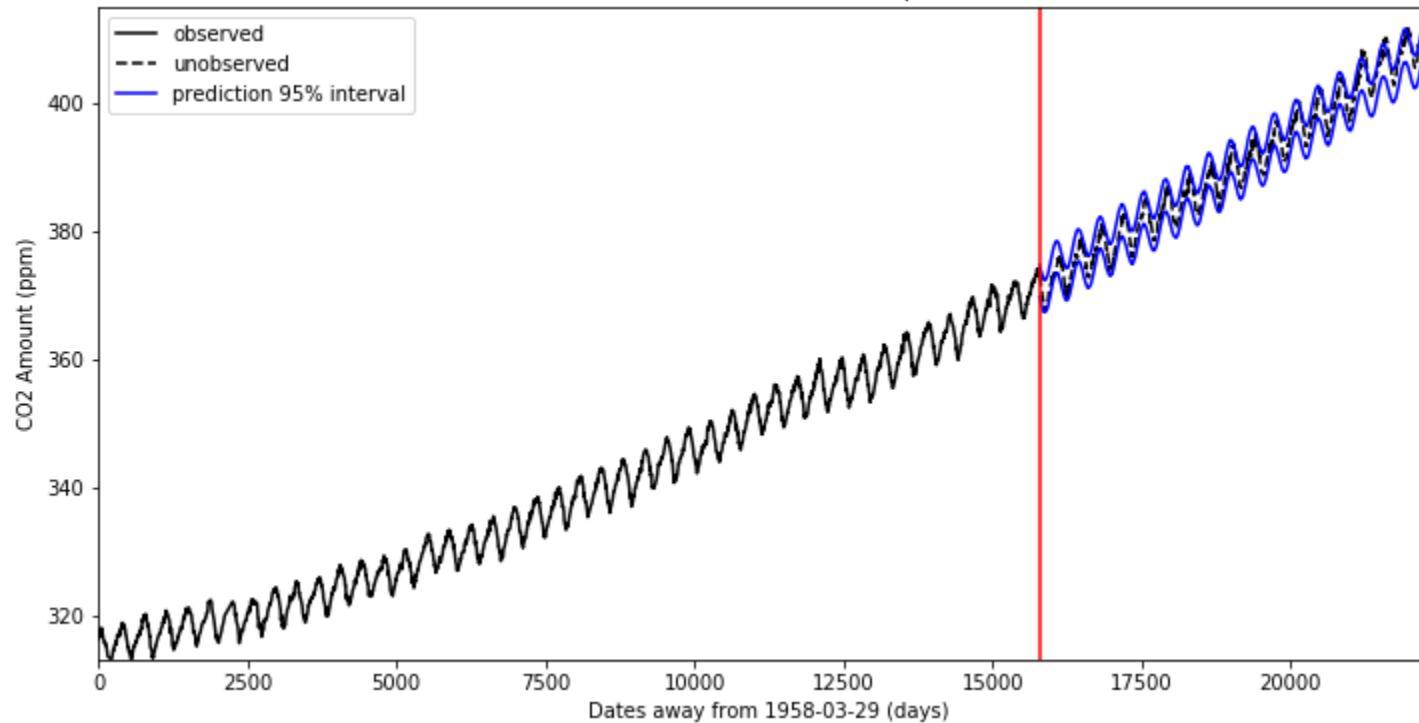
1 # draw predictions
2 prediction = samples['x_future']
3
4 # Compute 95% interval of the predicted values
5 prediction_interval = np.percentile(prediction, axis=0, q=[2.5, 97.5])
6
7 # Plot mean and 95% interval of predictions
8 plt.figure(figsize=(12, 6))
9 plt.plot(data["Time"][:n], data["Amount"][:n], 'k-', label='observed')
10 plt.plot(data["Time"][n:n+n_future], data["Amount"][n:n+n_future], 'k--', label='unobserved')
11 plt.xlim(0, data["Time"][n+n_future-1])
12 plt.ylim(
13     min(min(prediction[0,:]), min(data["Amount"])),
14     max(max(prediction[1,:]), max(data["Amount"])))
15 plt.plot(
16     data["Time"][n:],
17     prediction_interval[0,:],
18     'b-', label='prediction 95% interval')
19 plt.plot(

```

```
20     data["Time"][n:],
21     prediction_interval[1,:],
22     'b-')
23 plt.axvline(data["Time"][n - 1], color='red')
24 plt.xlabel("Dates away from 1958-03-29 (days)")
25 plt.ylabel("CO2 Amount (ppm)")
26 plt.legend()
27 plt.title('Observed data, unobserved future data, and predicted 95% interval')
28 plt.show()
```



Observed data, unobserved future data, and predicted 95% interval



```
1 # we store this sampling for later comparison
2 cos_prediction_interval = prediction_interval
3 cos_prediction = prediction
```

▼ New model: quadratic + sawtooth periodic

```
1 # The Stan model. Running this cell compiles the Stan model, which takes
2 # some time.
3 # Stan gets to see the first n data values only
4
5 stan_code = """
6 data {
7     int<lower=0> n;           // The number of data
8     int<lower=0> n_future;    // the number of predictions
9     real x[n];               // The data
10    real t[n];                // time data
11    real t_future[n_future]; // future time data
12 }
13
14 parameters {
15     real<lower=0> A;
16     real<lower=0> sigma;
17     real<lower = 0> a;
18     real<lower = 0> b;
19     real<lower = 0> c;
20     real phi_x;
21     real phi_y;
22 }
23
24 transformed parameters {
25     real<lower=0,upper=1> phi;
26     real<upper = 0> log_a;
27     // The atan2 function returns a value in the range [-pi, pi], which we
28     // then transform linearly into the range [0, 1].
29     phi = atan2(phi_x, phi_y) / (2 * pi()) + 0.5;
30
31     log_a = log(a);
32 }
33
34 model {
35     A ~ normal(0, 5);
36     sigma ~ gamma(1, 0.3);
37     log_a ~ normal(-16, 6);
38     b ~ normal(0, 0.1);
```

```

39     c ~ normal(300, 100);
40     phi_x ~ normal(0, 1);
41     phi_y ~ normal(0, 1);
42     for(i in 1:n) {
43         x[i] ~ normal(
44             A * (4 * fabs(fmod(1/365.25 * t[i] + phi, 1) - 0.5) - 1) + a*t[i]*t[i] + b*t[i] + c,
45             sigma);
46     }
47 }
48
49 // Generate the predicted function values for the next n_future steps.
50 generated quantities {
51     real x_future[n_future];
52     for(i in 1:n_future) {
53         x_future[i] = normal_rng(
54             A * (4 * fabs(fmod(1/365.25 * t_future[i] + phi, 1) - 0.5) - 1) + a*t_future[i]*t_future[i] + b*t_future[i] + c,
55             sigma);
56     }
57 }
58 """
59
60 stan_model = pystan.StanModel(model_code=stan_code)

```

↳ INFO:pystan:COMPIILING THE C++ CODE FOR MODEL anon_model_e495cf1d5c27bdec1a1bdddf6401e82a NOW.

```

1 # subsetting the data
2 n = len(data)*7//10
3 n_future = len(data) - n
4
5 stan_data = {
6     'n': n,
7     'n_future': n_future,
8     'x': data["Amount"][:n],
9     't': data["Time"][:n],
10    "t_future" : data["Time"][n:]}
11
12 # Run Hamiltonian Monte Carlo using Stan with 4 Markov chains, a 1000 step
13 # warm-up phase and 1000 step sampling phase for each chain. The warm-up samples
14 # are discarded, so we are left with 4 x 1000 = 4000 samples.

```

```

15
16 # This cell will take a minute or two to run.
17
18 parameters = ['A', 'phi', 'a', "b", "c", "sigma"]
19
20 results = stan_model.sampling(data=stan_data)
21 print(results.stansummary(pars=parameters))
22 samples = results.extract()

```

↳ WARNING:pystan:73 of 4000 iterations ended with a divergence (1.82 %).
 WARNING:pystan:Try running with adapt_delta larger than 0.8 to remove the divergences.
 Inference for Stan model: anon_model_e495cf1d5c27bdec1a1bdddf6401e82a.
 4 chains, each with iter=2000; warmup=1000; thin=1;
 post-warmup draws per chain=1000, total post-warmup draws=4000.

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
A	3.41	8.1e-4	0.04	3.34	3.39	3.41	3.43	3.48	1997	1.0
phi	0.93	2.8e-5	1.7e-3	0.93	0.93	0.93	0.93	0.94	3661	1.0
a	9.0e-8	3.1e-11	1.1e-9	8.7e-8	8.9e-8	9.0e-8	9.1e-8	9.2e-8	1355	1.0
b	2.2e-3	5.1e-7	1.9e-5	2.2e-3	2.2e-3	2.2e-3	2.2e-3	2.3e-3	1375	1.0
c	314.2	1.6e-3	0.07	314.07	314.16	314.2	314.25	314.33	1694	1.0
sigma	0.98	3.3e-4	0.01	0.95	0.97	0.98	0.99	1.01	2011	1.0

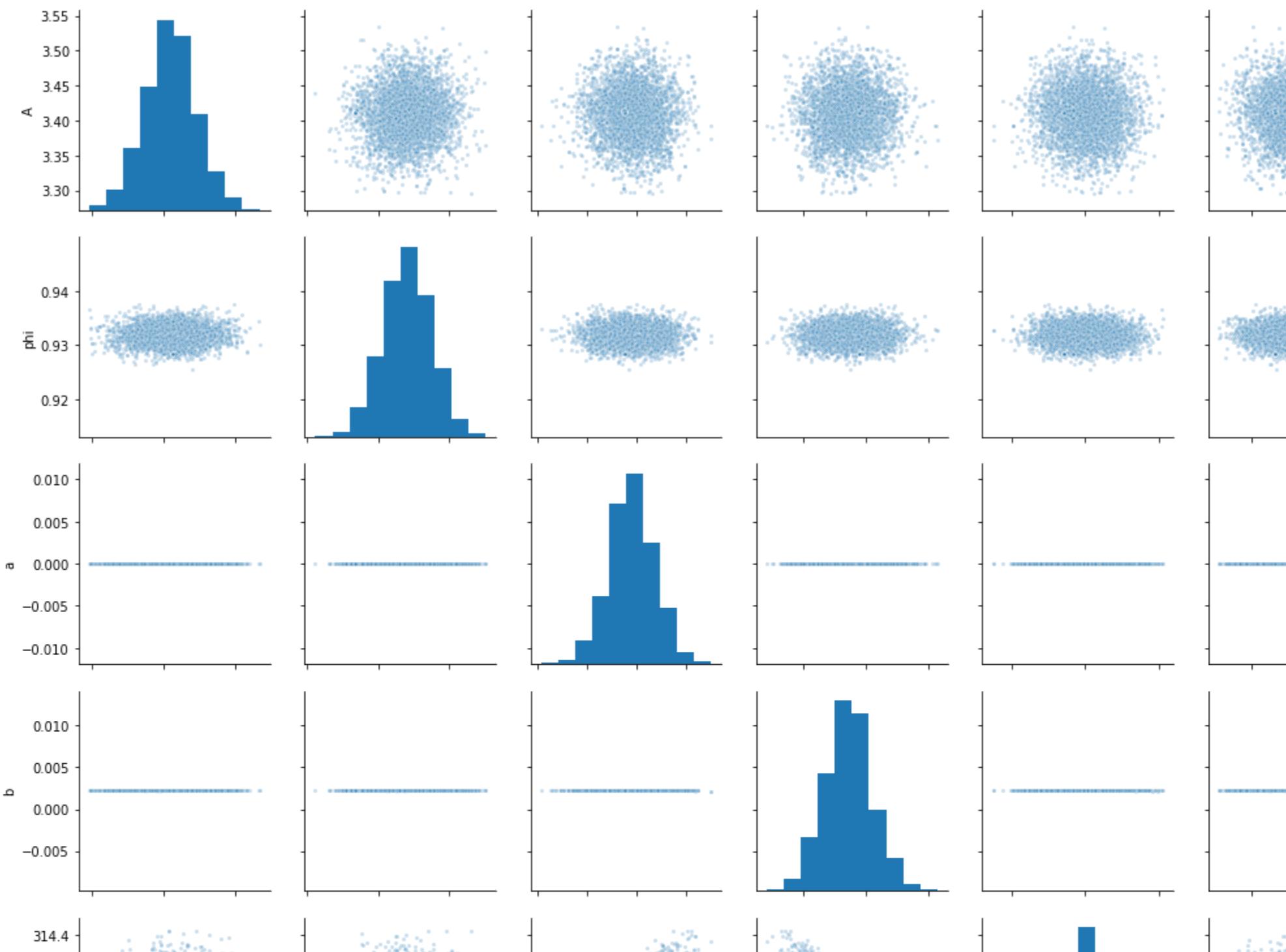
Samples were drawn using NUTS at Wed Dec 18 14:12:00 2019.
 For each parameter, n_eff is a crude measure of effective sample size,
 and Rhat is the potential scale reduction factor on split chains (at
 convergence, Rhat=1).

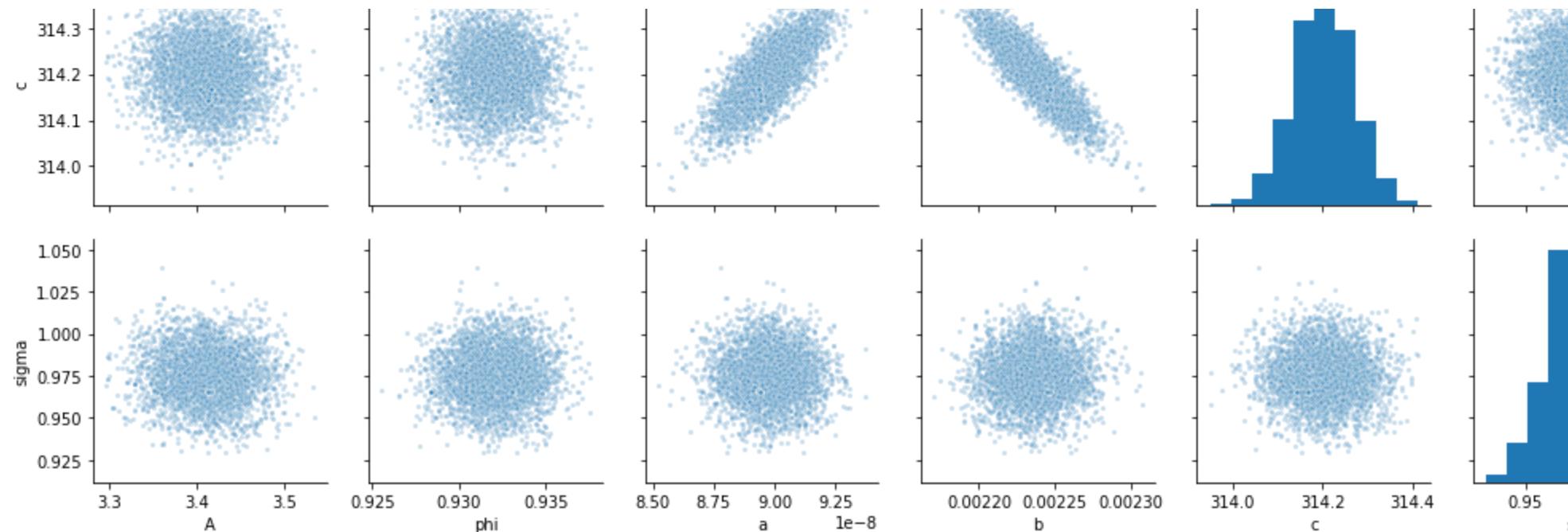
```

1 df = pd.DataFrame(
2     data=np.transpose([samples[param] for param in parameters]),
3     columns=parameters)
4 seaborn.pairplot(df, height=2.5, plot_kws={'marker': '.', 'alpha': 0.25})
5 plt.show()

```

↳





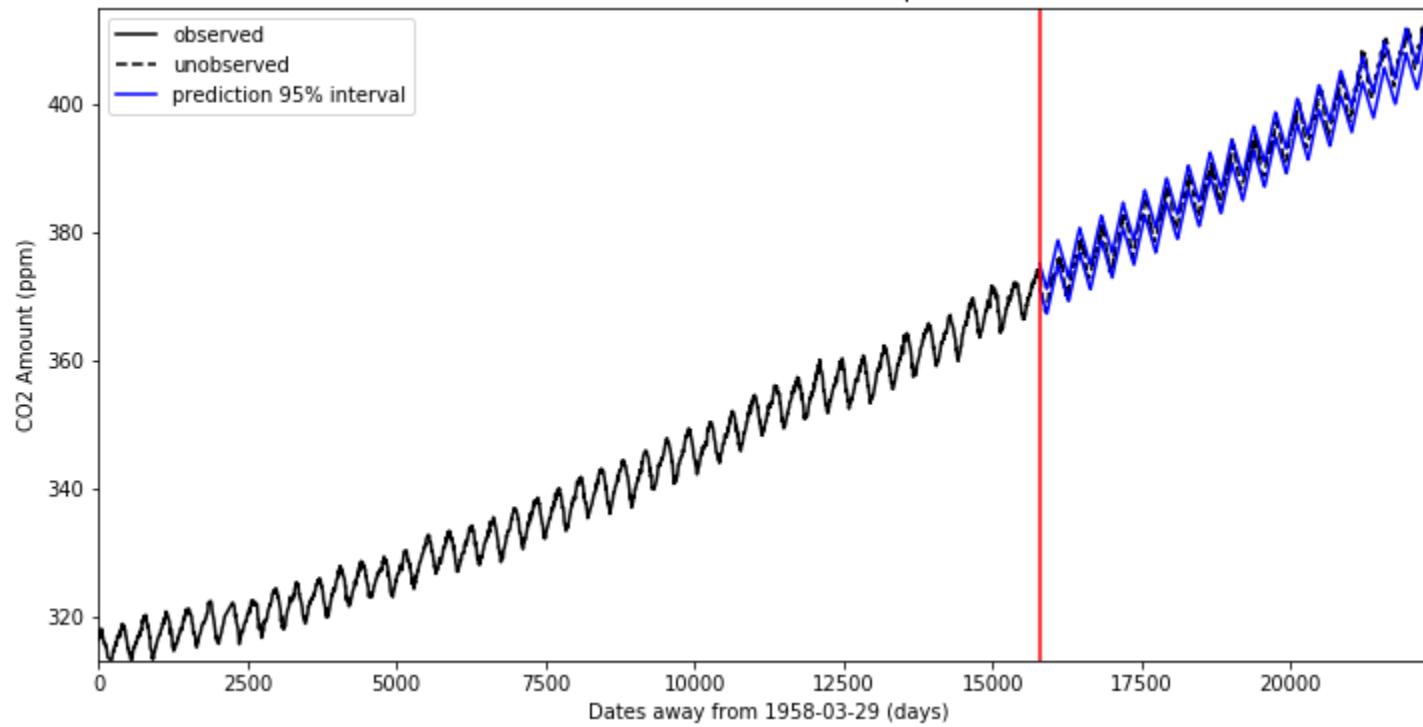
```

1 prediction = samples['x_future']
2
3 # Compute 95% interval of the predicted values
4 prediction_interval = np.percentile(prediction, axis=0, q=[2.5, 97.5])
5
6 # Plot mean and 95% interval of predictions
7 plt.figure(figsize=(12, 6))
8 plt.plot(data["Time"][:n], data["Amount"][:n], 'k-', label='observed')
9 plt.plot(data["Time"][n:n+n_future], data["Amount"][n:n+n_future], 'k--', label='unobserved')
10 plt.xlim(0, data["Time"][n+n_future-1])
11
12 # set scaling
13 plt.ylim(
14     min(min(prediction[0,:]), min(data["Amount"])),
15     max(max(prediction[1,:]), max(data["Amount"])))
16
17 # plot confidence interval
18 plt.plot(
19     data["Time"][n:],
```

```
20     prediction_interval[0,:],  
21     'b-', label='prediction 95% interval')  
22 plt.plot(  
23     data["Time"][n:],  
24     prediction_interval[1,:],  
25     'b-')  
26  
27 plt.axvline(data["Time"][n - 1], color='red')  
28 plt.xlabel("Dates away from 1958-03-29 (days)")  
29 plt.ylabel("CO2 Amount (ppm)")  
30 plt.legend()  
31 plt.title('Observed data, unobserved future data, and predicted 95% interval')  
32 plt.show()
```



Observed data, unobserved future data, and predicted 95% interval



```
1 # we store this sampling for later comparison  
2 saw_prediction_interval = prediction_interval  
3 saw_prediction = prediction
```

▼ Model Comparison

```
1 # confidence interval width comparison
2 saw_interval_length = saw_prediction_interval[1:,:][0] - saw_prediction_interval[0:,:][0]
3 cos_interval_length = cos_prediction_interval[1:,:][0] - cos_prediction_interval[0:,:][0]
4 rate = np.mean(saw_interval_length < cos_interval_length)*100
5 print("The percentage of the width of model 3 interval is smaller than of model 2 interval: {}%".format(rate))
```

⇒ The percentage of the width of model 3 interval is smaller than of model 2 interval: 100.0%

▼ Accuracy testing

```
1 # the real data to validate
2 validate = np.array(data["Amount"][:n])
3
4 # the mean of prediction
5 saw_mean = np.mean(saw_prediction, axis = 0)
6 cos_mean = np.mean(cos_prediction, axis = 0)
7
8 # compute the SME errors
9 saw_SME = np.mean((saw_mean - validate)**2)
10 cos_SME = np.mean((cos_mean - validate)**2)
11 saw_all_SME = np.mean((saw_prediction - validate)**2)
12 cos_all_SME = np.mean((cos_prediction - validate)**2)
13
14 print("Sawtooth SME wrt to samples mean:", saw_SME)
15 print("Cos SME wrt to samples mean:", cos_SME)
16 print("Sawtooth SME wrt to all samples values:", saw_all_SME)
17 print("Cos SME wrt to all samples values:", cos_all_SME)
```

⇒ Sawtooth SME wrt to samples mean: 1.612156603044492
Cos SME wrt to samples mean: 2.1612224291754356
Sawtooth SME wrt to all samples values: 2.5857090803350777
Cos SME wrt to all samples values: 3.824635218251834

▼ Test statistics

As discussed above, we use the first 70% of the data (2197 values) as data for computing the posterior distributions over the parameters above using Stan. We use the other 30% of the data (the next 942 values) to test our prediction of future function values.

Our test statistic is the mean of future data: 942 last data points in the dataset.

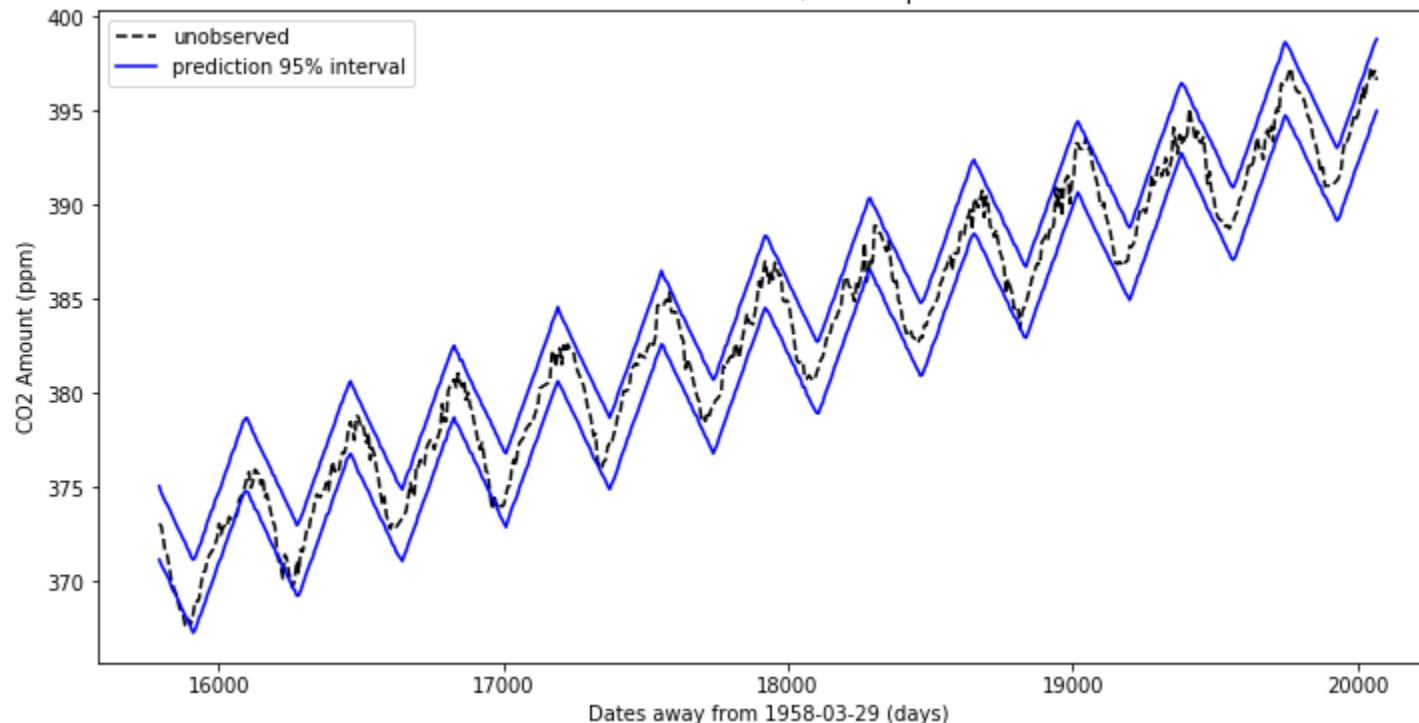
The reason is that: we suspect the model, even though overall fit the future data very well (based on the plot), fails to predict the few anomalies in the original data fully. For example, when we zoom in the first 600 data points in the test future data, we can see at those peaks, our real data is slightly off compared to the predicted interval.

As the mean is sensitive to the outliers, we would expect the mean for the real future data to be in the upper or bottom percentile of the test statistics results of the sampled data.

```
1 # zoom in the data for the first 600 points
2 plt.figure(figsize = (12,6))
3 plt.plot(data["Time"][n:n+600], data["Amount"][n:n+600], 'k--', label='unobserved')
4 plt.plot(
5     data["Time"][n:n+600],
6     prediction_interval[0,:][:600],
7     'b-', label='prediction 95% interval')
8 plt.plot(
9     data["Time"][n:n+600],
10    prediction_interval[1,:][:600],
11    'b-')
12 plt.xlabel("Dates away from 1958-03-29 (days)")
13 plt.ylabel("CO2 Amount (ppm)")
14 plt.legend()
15 plt.title('The first 600 unobserved future data, and its predicted 95% interval')
16 plt.show()
```



The first 600 unobserved future data, and its predicted 95% interval



```

1 # test statistics: mean of samples
2 mean_samples = np.mean(saw_prediction, axis = 1)
3 benchmark = np.mean(data["Amount"][:n])
4
5 # plot the histogram of the test statistics
6 plt.figure(figsize = (12,6))
7 plt.hist(mean_samples, alpha = 0.5)
8 plt.axvline(np.quantile(mean_samples, 0.025), color='red', label = "95% interval", linestyle = "--")
9 plt.axvline(np.quantile(mean_samples, 0.975), color='red', linestyle = "--")
10 plt.axvline(np.mean(mean_samples), color='blue', linestyle = "--", label = "Mean samples")
11 plt.axvline(benchmark, color='red', label = "Real data = " + str(benchmark))
12 plt.xlabel("CO2 Amount (ppm)")
13 plt.title("Histogram of the test statistics of the simulated data from model 3\n compared to\n the test statistics for the real data (red line).")
14 plt.legend()
15 plt.show()

```

17

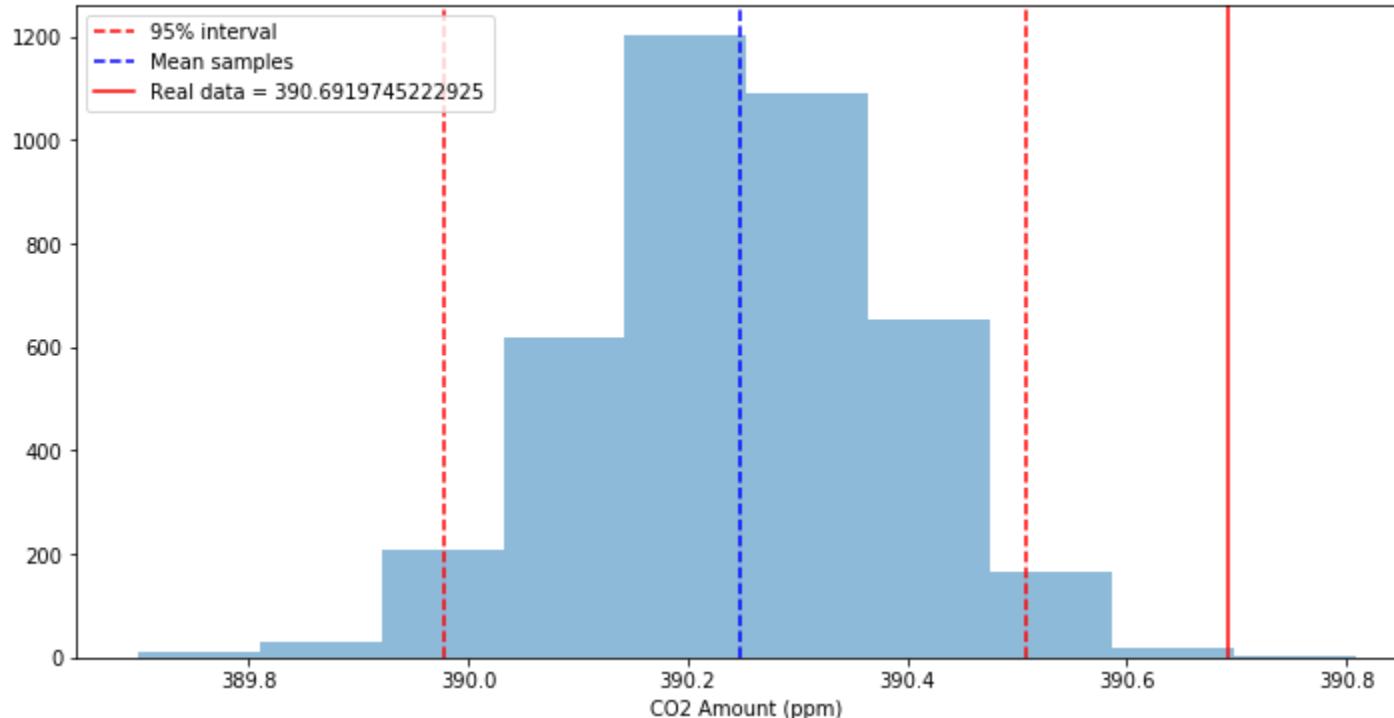
```

18 final = np.mean(mean_samples > benchmark)
19 print("percentage result > real statistic value: ", final*100, "%")
20 print("p value: ", 1-final)

```

C

Histogram of the test statistics of the simulated data from model 3 compared to the test statistics for the real data (red line).



percentage result > real statistic value: 0.1 %
p value: 0.999

```

1 # test statistics: mean of samples
2 mean_samples = np.mean(cos_prediction, axis = 1)
3 benchmark = np.mean(data["Amount"][:n])
4
5 # plot the histogram of the test statistics
6 plt.figure(figsize = (12,6))
7 plt.hist(mean_samples, alpha = 0.5)
8 plt.axvline(np.quantile(mean_samples, 0.025), color='red', label ="95% interval", linestyle = "--")
9 plt.axvline(np.quantile(mean_samples, 0.975), color='red', linestyle = "--")

```

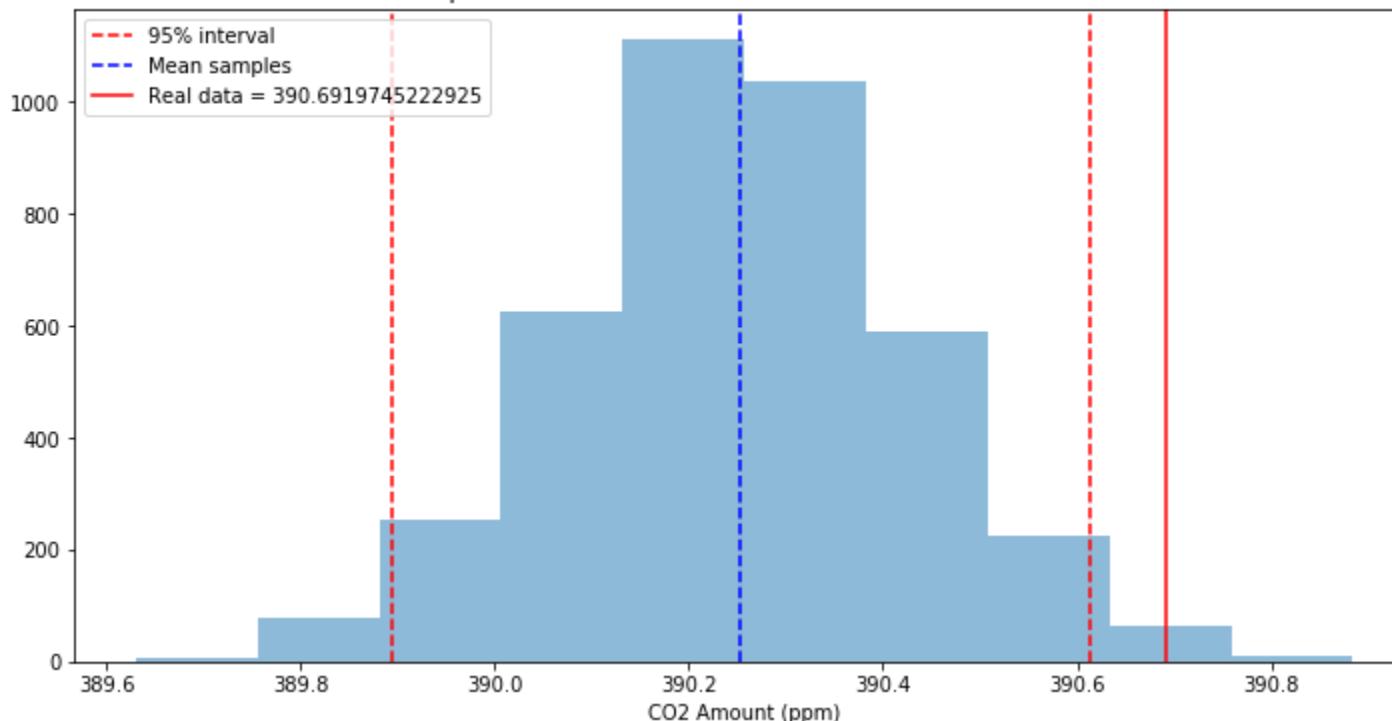
```

10 plt.axvline(np.mean(mean_samples), color='blue', linestyle = "--", label = "Mean samples")
11 plt.axvline(benchmark, color='red', label ="Real data = " + str(benchmark))
12 plt.xlabel("CO2 Amount (ppm)")
13 plt.title("Histogram of the test statistics of the simulated data from model 2\n compared to\
14 the test statistics for the real data (red line).")
15 plt.legend()
16 plt.show()
17
18 final = np.mean(mean_samples > benchmark)
19 print("percentage result > real statistic value: ", final*100, "%")
20 print("p value: ", 1-final)

```



Histogram of the test statistics of the simulated data from model 2
compared to the test statistics for the real data (red line).



percentage result > real statistic value: 0.7250000000000001 %
p value: 0.99275

```

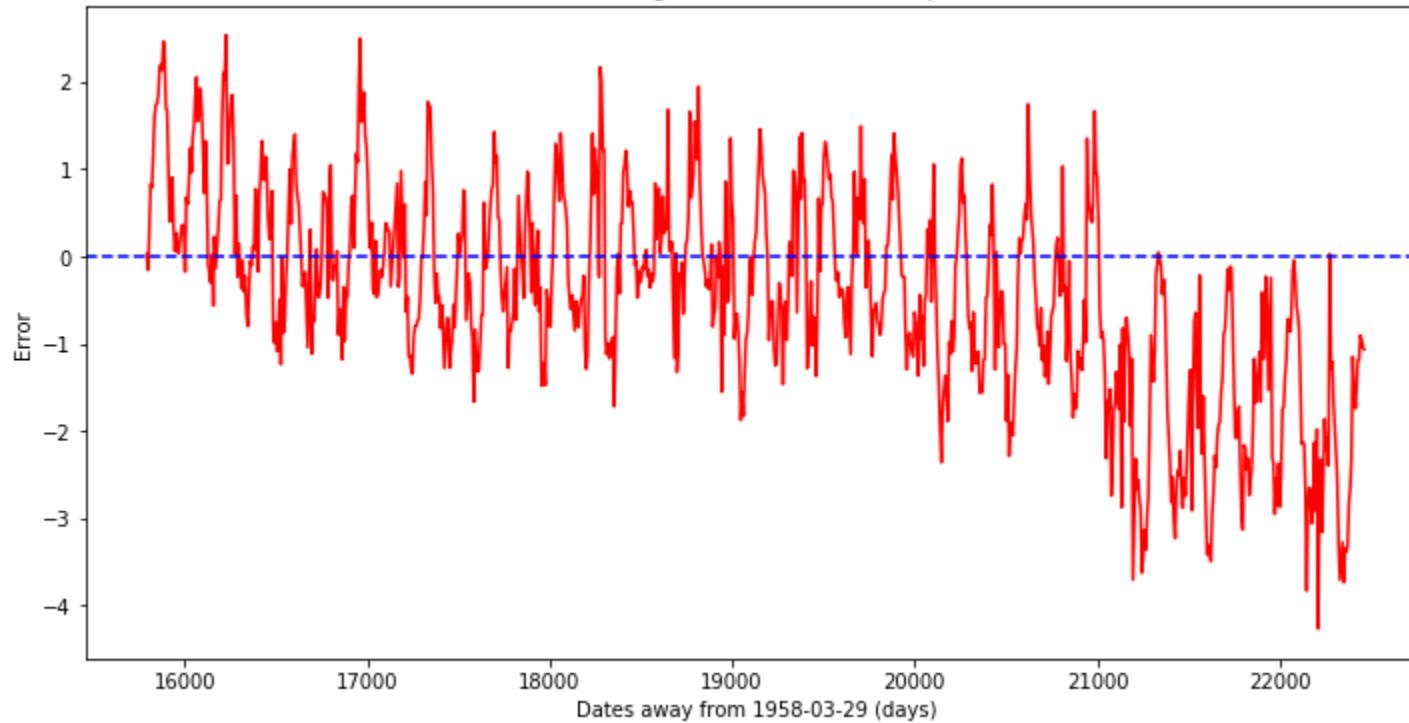
1 # the mean of errors by time
2 error_by_time = np.mean(saw_prediction - validate, axis = 0)
3

```

```
4 # plot
5 plt.figure(figsize = (12,6))
6 plt.plot(data["Time"][n:], error_by_time, 'r-', label = "error")
7 plt.axhline(0, color = "blue", linestyle = "--")
8 plt.xlabel("Dates away from 1958-03-29 (days)")
9 plt.ylabel("Error")
10 plt.title('Errors throughout the unobserved period')
11 plt.show()
```



Errors throughout the unobserved period



```
1 # combine of every errors
2 error_sets = (saw_prediction - validate).flatten()
3 print("Median:", np.median(error_sets))
4 print("95% Interval:", np.quantile(error_sets, [0.025, 0.975]))
5
6 # plot
7 plt.figure(figsize = (12,6))
8 plt.hist(error_sets)
```

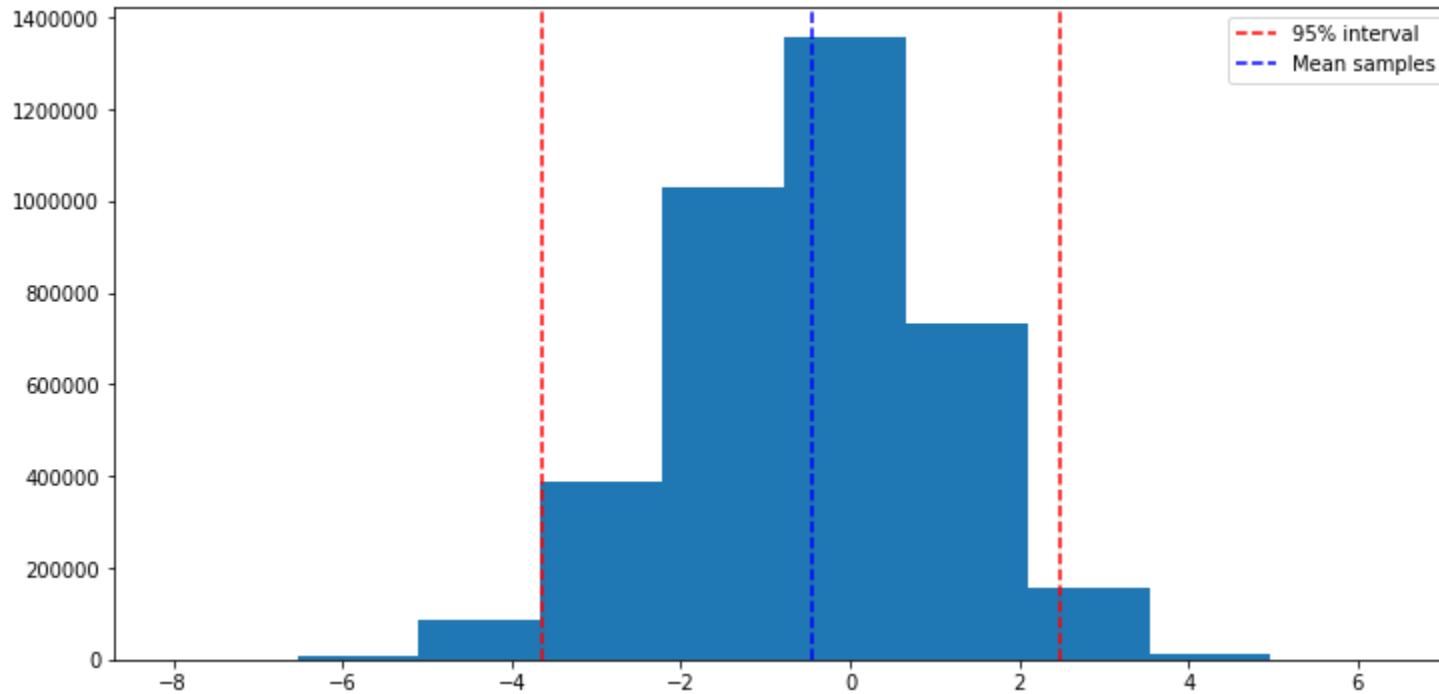
```

9 plt.axvline(np.quantile(error_sets, 0.025), color='red', label ="95% interval", linestyle = "--")
10 plt.axvline(np.quantile(error_sets, 0.975), color='red', linestyle = "--")
11 plt.axvline(np.mean(error_sets), color='blue', linestyle = "--", label = "Mean samples")
12 plt.legend()
13 plt.show()

```

Median: -0.39648976903481525

95% Interval: [-3.63786624 2.4763771]



▼ Predictions

```

1 # the dates need for calculation: 1st date of data, end date of data, 2058 dates (start and end)
2 first_date = datetime.datetime.strptime(data["Date"][0], "%Y-%m-%d")
3 last_date = datetime.datetime.strptime(data["Date"][-1], "%Y-%m-%d")
4 desire_first = datetime.datetime.strptime("2058-01-01", "%Y-%m-%d")
5 desire_last = datetime.datetime.strptime("2058-12-31", "%Y-%m-%d")
6
7 # convert into readable dates by pyStan
8 start1 = (desire_first - first_date).days // 7 + 1

```

```
0 stop1 = (desire_end - desire_start).days//7
1 stop2 = (desire_last - first_date).days//7
2 start = (last_date - first_date).days//7
3
4
5 # generate the future time parameters
6 new_time = np.array(list(range(int(start + 1), int(start + 1) + int(40*365.25//7))))*7
7 new_time = list(new_time)

8 # rerun the model with all the data to predict the future
9 n = len(data)
10 n_future = len(new_time)
11
12 stan_data = {
13     'n': n,
14     'n_future': n_future,
15     'x': data["Amount"],
16     't': data["Time"],
17     "t_future" : new_time}
18
19
20 # Run Hamiltonian Monte Carlo using Stan with 4 Markov chains, a 1000 step
21 # warm-up phase and 1000 step sampling phase for each chain. The warm-up samples
22 # are discarded, so we are left with 4 x 1000 = 4000 samples.
23
24 # This cell will take a minute or two to run.
25
26 parameters = ['A', 'phi', 'log_a', 'b', 'c', "sigma"]
27
28 results = stan_model.sampling(data=stan_data)
29 print(results.stansummary(pars=parameters))
30 samples = results.extract()
```



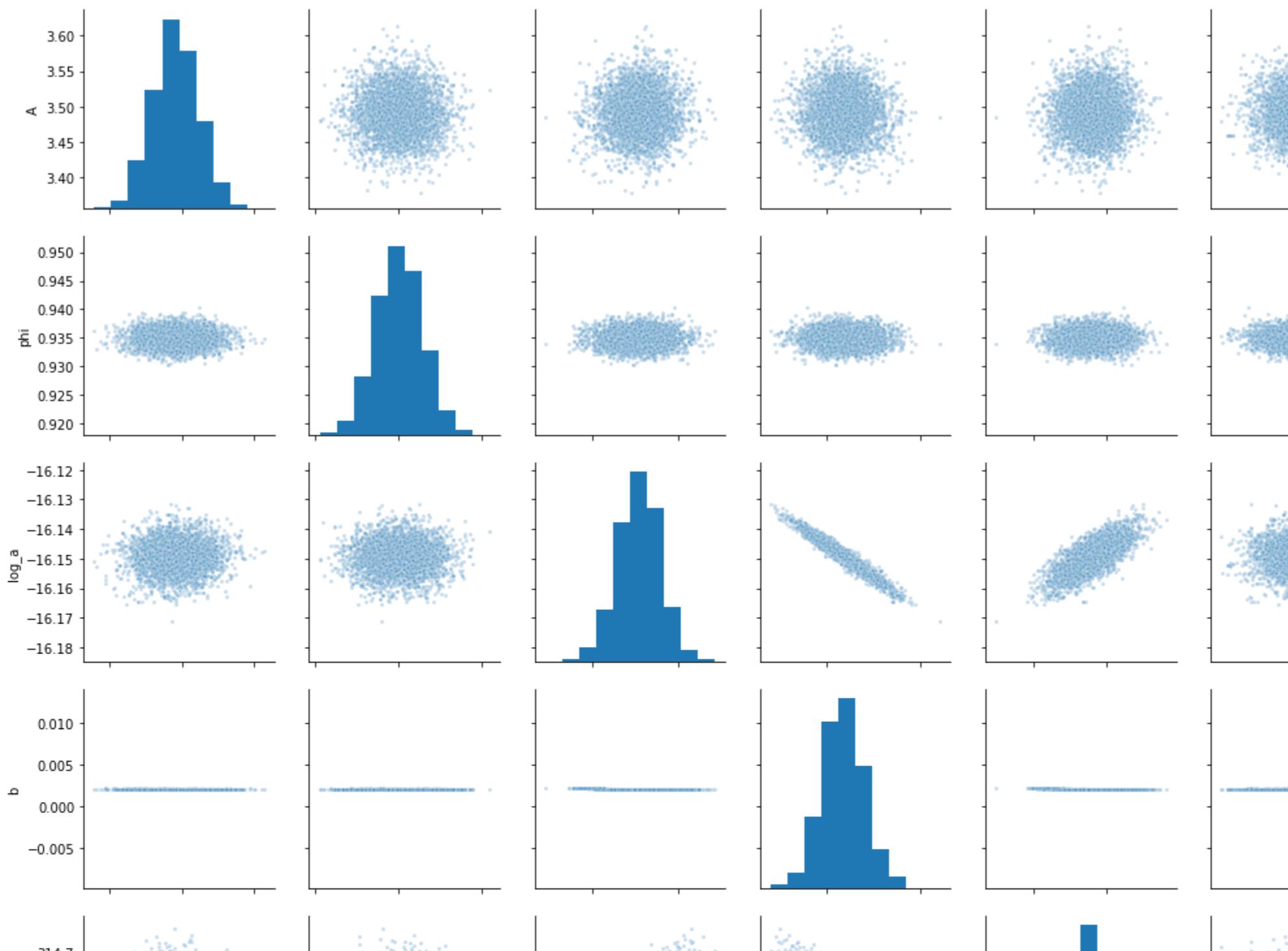
WARNING:pystan:Maximum (flat) parameter count (1000) exceeded: skipping diagnostic tests for n_eff and Rhat.
To run all diagnostics call pystan.check_hmc_diagnostics(fit)
WARNING:pystan:54 of 4000 iterations ended with a divergence (1.35 %).
WARNING:pystan:Try running with adapt_delta larger than 0.8 to remove the divergences.
WARNING:pystan:100 of 4000 iterations saturated the maximum tree depth of 10 (2.5 %)
WARNING:pystan:Run again with max_treedepth larger than 10 to avoid saturation
Inference for Stan model: anon_model_e495cf1d5c27bdec1a1bdddf6401e82a.
4 chains, each with iter=2000; warmup=1000; thin=1;
post-warmup draws per chain=1000, total post-warmup draws=4000.

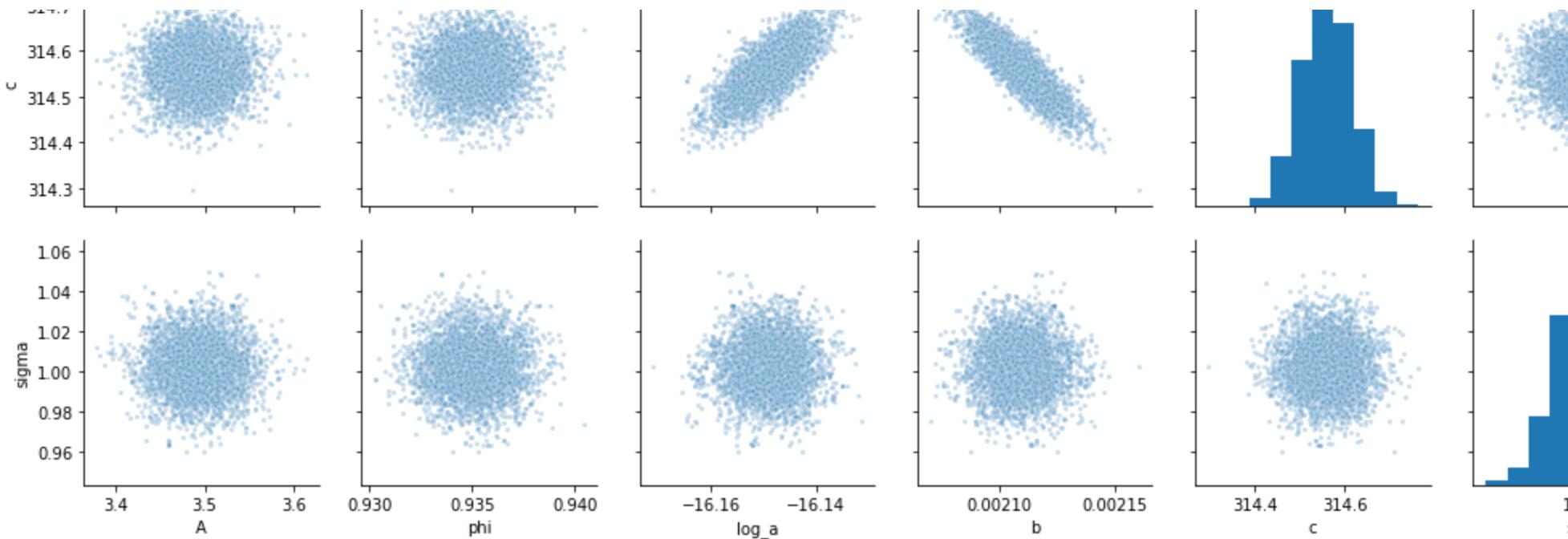
	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
A	3.49	7.2e-4	0.03	3.43	3.47	3.49	3.51	3.55	1939	1.0
phi	0.94	2.2e-5	1.4e-3	0.93	0.93	0.94	0.94	0.94	4093	1.0
log_a	-16.15	1.2e-4	5.1e-3	-16.16	-16.15	-16.15	-16.15	-16.14	1665	1.0
b	2.1e-3	2.8e-7	1.1e-5	2.1e-3	2.1e-3	2.1e-3	2.1e-3	2.1e-3	1696	1.0
c	314.56	1.2e-3	0.06	314.45	314.52	314.56	314.6	314.67	2152	1.0
sigma	1.0	3.0e-4	0.01	0.98	0.99	1.0	1.01	1.03	1881	1.0

Samples were drawn using NUTS at Wed Dec 18 14:32:04 2019.
For each parameter, n_eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor on split chains (at
convergence, Rhat=1).

```
1 # pair plots
2 df = pd.DataFrame(
3     data=np.transpose([samples[param] for param in parameters]),
4     columns=parameters)
5 seaborn.pairplot(df, height=2.5, plot_kws={'marker': '.', 'alpha': 0.25})
6 plt.show()
```







```

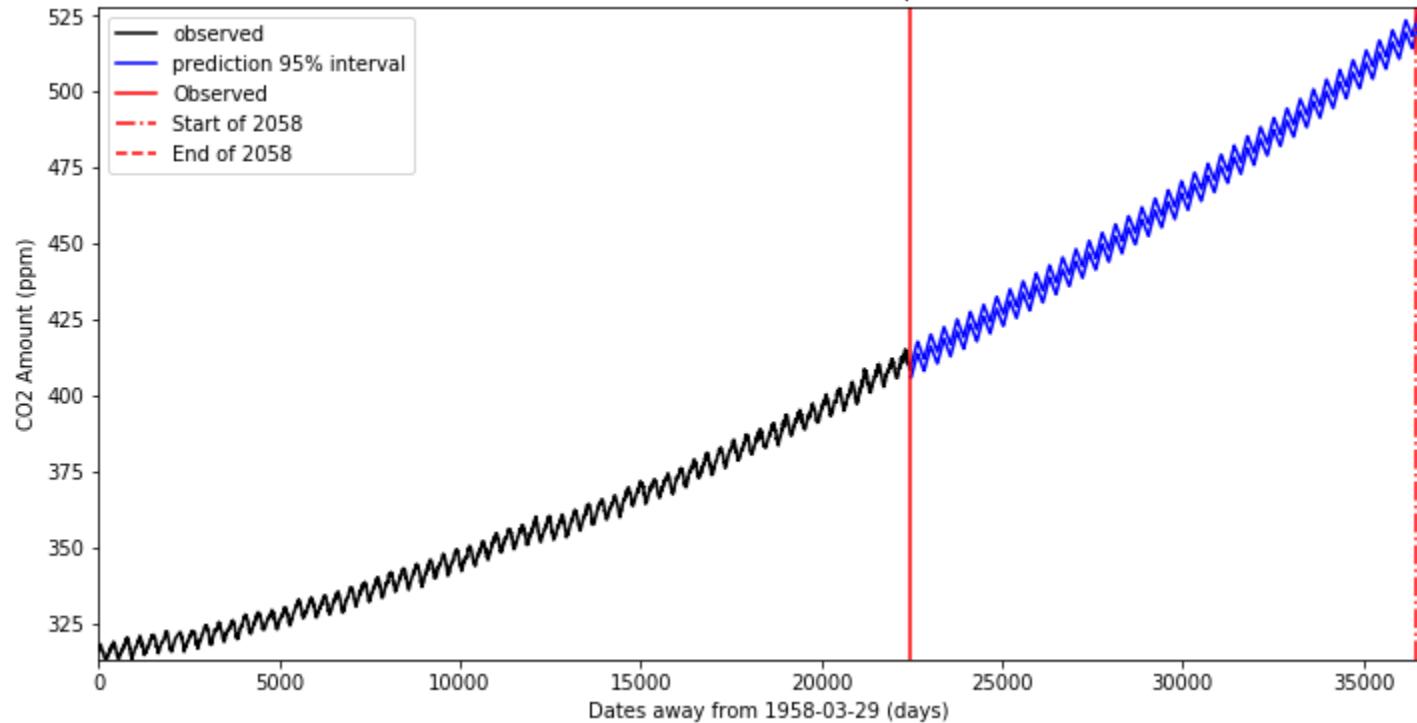
1 prediction = samples['x_future']
2
3 # Compute 95% interval of the predicted values
4 prediction_interval = np.percentile(prediction, axis=0, q=[2.5, 97.5])
5
6 # Plot mean and 95% interval of predictions
7 plt.figure(figsize=(12, 6))
8 plt.plot(data["Time"], data["Amount"], 'k-', label='observed')
9 plt.xlim(0, max(new_time))
10 plt.ylim(
11     min(min(prediction[0,:]), min(data["Amount"])),
12     max(max(prediction[1,:]), max(data["Amount"])))
13
14 plt.plot(
15     new_time,
16     prediction_interval[0,:],
17     'b-', label='prediction 95% interval')
18
19 plt.plot(

```

```
20     new_time,
21     prediction_interval[1,:],
22     'b-')
23
24 plt.axvline(data["Time"][n - 1], color='red', label = "Observed")
25 plt.axvline(stop1 * 7, color = "red", linestyle = "-.", label ="Start of 2058")
26 plt.axvline(stop2 * 7, color = "red", linestyle = "--", label ="End of 2058")
27 plt.xlabel("Dates away from 1958-03-29 (days)")
28 plt.ylabel("CO2 Amount (ppm)")
29 plt.legend()
30 plt.title('Observed data, unobserved future data, and predicted 95% interval')
31 plt.show()
```



Observed data, unobserved future data, and predicted 95% interval



▼ year 2058

```
1 # Get data for this specific measurement and compute the statistics
```

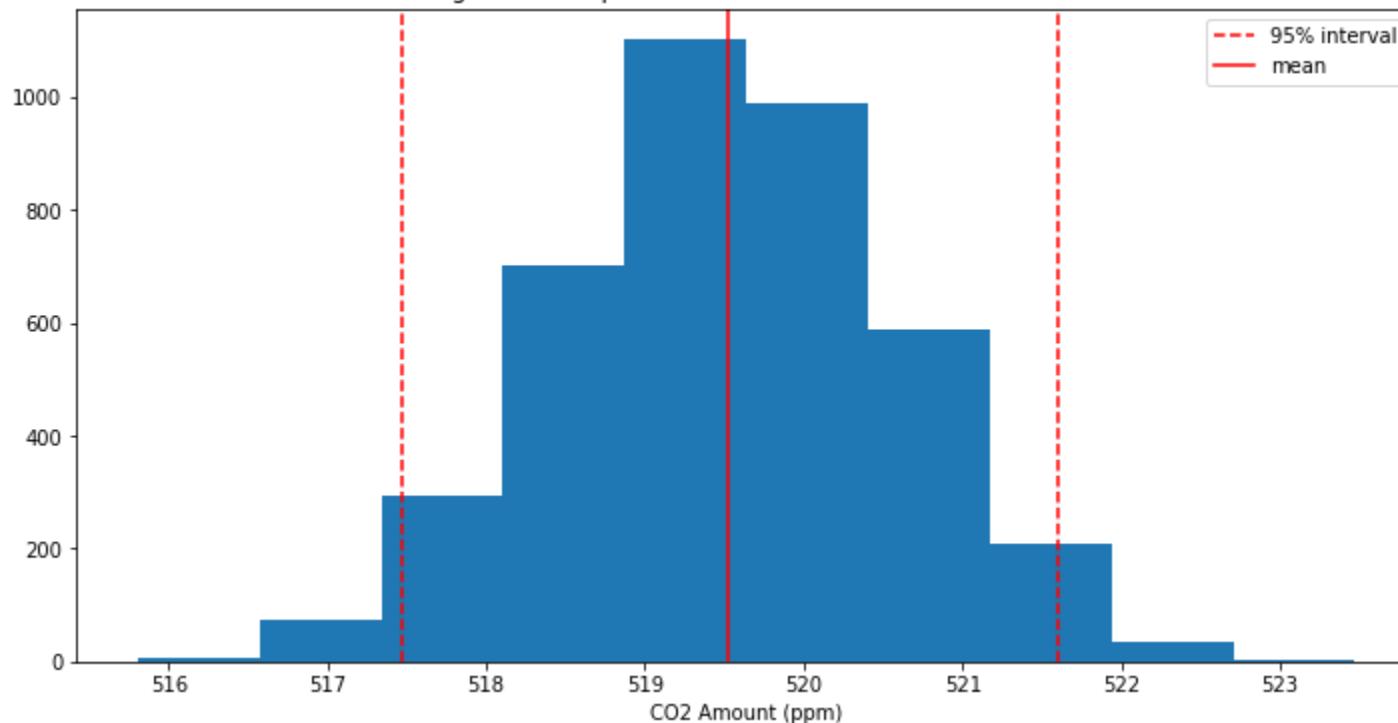
<https://colab.research.google.com/drive/1KWIUWOi9tH63n6MPaxxKX8lnGOQCMSk?authuser=1#scrollTo=KGO8Kw2-v8DS&printMode=true>

```
2 sample_2058 = prediction.T[int(stop1-start)-1]
3 print("Mean of the estimate: {}".format(np.mean(sample_2058)))
4 print("Median of the estimate: {}".format(np.median(sample_2058)))
5 print("STD of the estimate:", np.std(sample_2058))
6 lower_2058 = prediction_interval[0,:][int(stop1-start)-1]
7 upper_2058 = prediction_interval[1,:][int(stop1-start)-1]
8 print("Confidence interval in 2058: [{}, {}]".format(lower_2058, upper_2058))
9
10 # plot the results
11 plt.figure(figsize = (12,6))
12 plt.hist(sample_2058)
13 plt.xlabel("CO2 Amount (ppm)")
14 plt.title("Histogram of the predicted values for CO2 Amount in 2058")
15 plt.axvline(lower_2058, color = "red", linestyle = "--")
16 plt.axvline(upper_2058, color = "red", linestyle = "--", label = "95% interval")
17 plt.axvline(np.mean(sample_2058), color = "red", linestyle = "-", label = "mean")
18 plt.legend()
19 plt.show()
```



Mean of the estimate: 519.5246199645143
Median of the estimate: 519.5167440251244
STD of the estimate: 1.05966781226145
Confidence interval in 2058: [517.474873011385, 521.5981563346509]

Histogram of the predicted values for CO2 Amount in 2058



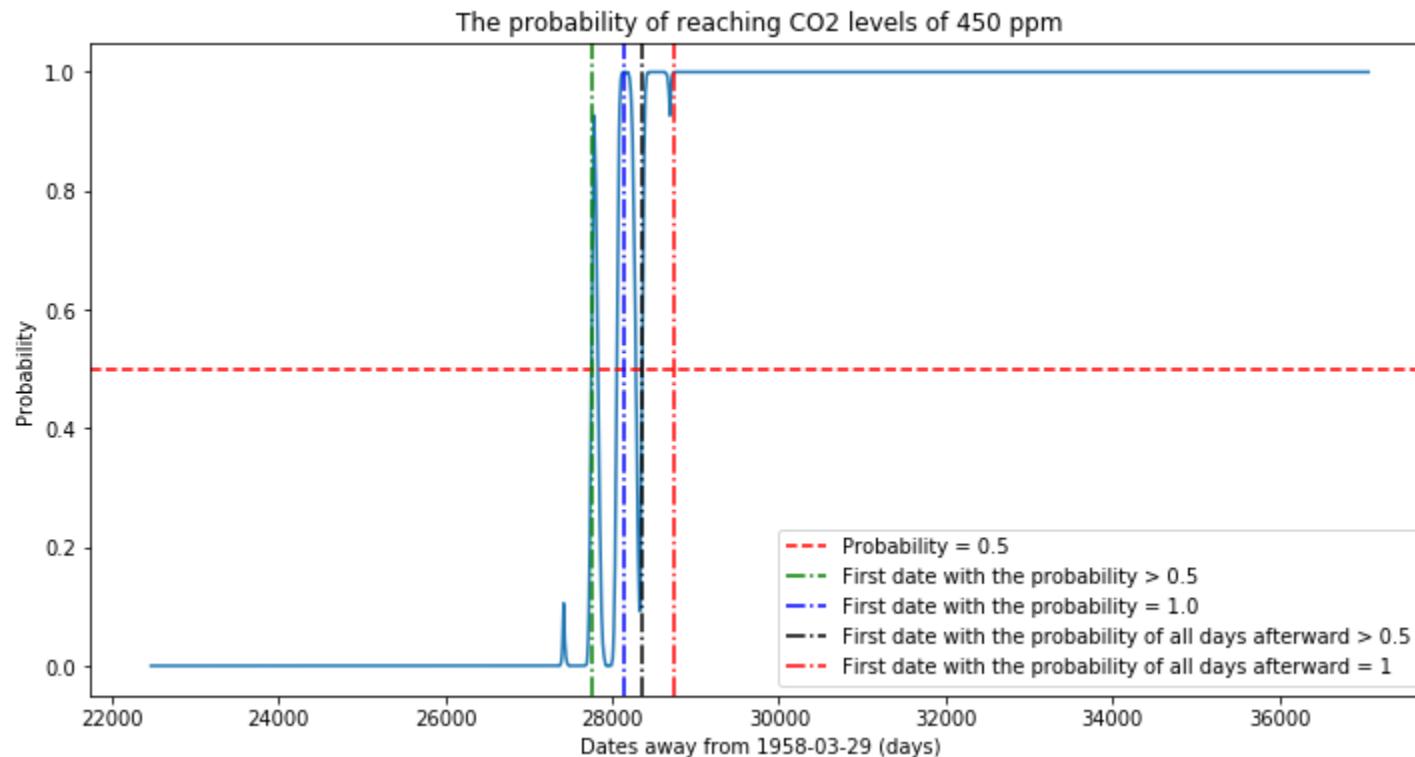
```
1 # the probabilities of certain dates > 450 ppm
2 probability = np.mean(prediction.T >= 450, axis = 1)
3
4 # first value to have 1.0 probability
5 idx_100 = next(x[0] for x in enumerate(list(probability)) if x[1] >= 1)
6
7 # first value to have 0.5 probability
8 idx_50 = next(x[0] for x in enumerate(list(probability)) if x[1] >= 0.5)
9
10 # first value to have all values onward = 1.0
11 for i in range(len(probability)):
12     if np.mean(probability[i:]) == 1:
```

```
13     idx_full_100 = i
14     break
15
16 # first value to have all values onward >= 0.5
17 for i in range(len(probability)):
18     if np.mean(probability[i:]) >= 0.5) == 1:
19         idx_full_50 = i
20         break
21
22 print("1st date to have 0.5 probability:", timedelta(idx_50*7) + last_date)
23 print("1st date to have 1.0 probability:", timedelta(idx_100*7) + last_date)
24 print("1st date to have all 0.5 probability onward:", timedelta(idx_full_50*7) + last_date)
25 print("1st date to have all 1.0 probability onward:", timedelta(idx_full_100*7) + last_date)
```

```
↳ 1st date to have 0.5 probability: 2034-03-18 00:00:00
    1st date to have 1.0 probability: 2035-04-07 00:00:00
    1st date to have all 0.5 probability onward: 2035-11-17 00:00:00
    1st date to have all 1.0 probability onward: 2036-11-29 00:00:00
```

```
1 # plot the probability
2 plt.figure(figsize = (12,6))
3 plt.plot(new_time,probability)
4 plt.xlabel("Dates away from 1958-03-29 (days)")
5 plt.ylabel("Probability")
6 plt.axhline(0.5, color = "red", linestyle = "--", label = "Probability = 0.5")
7 plt.axvline(idx_50 * 7 + data["Time"][n - 1], color = "green", linestyle = "-.", label ="First date with the probability > 0.5")
8 plt.axvline(idx_100 * 7 + data["Time"][n - 1], color = "blue", linestyle = "-.", label ="First date with the probability = 1.0")
9 plt.axvline(idx_full_50 * 7 + data["Time"][n - 1], color = "black", linestyle = "-.", label ="First date with the probability of all")
10 plt.axvline(idx_full_100 * 7 + data["Time"][n - 1], color = "red", linestyle = "-.", label ="First date with the probability of all")
11 plt.title("The probability of reaching CO2 levels of 450 ppm")
12 plt.legend()
13 plt.show()
```

```
↳
```

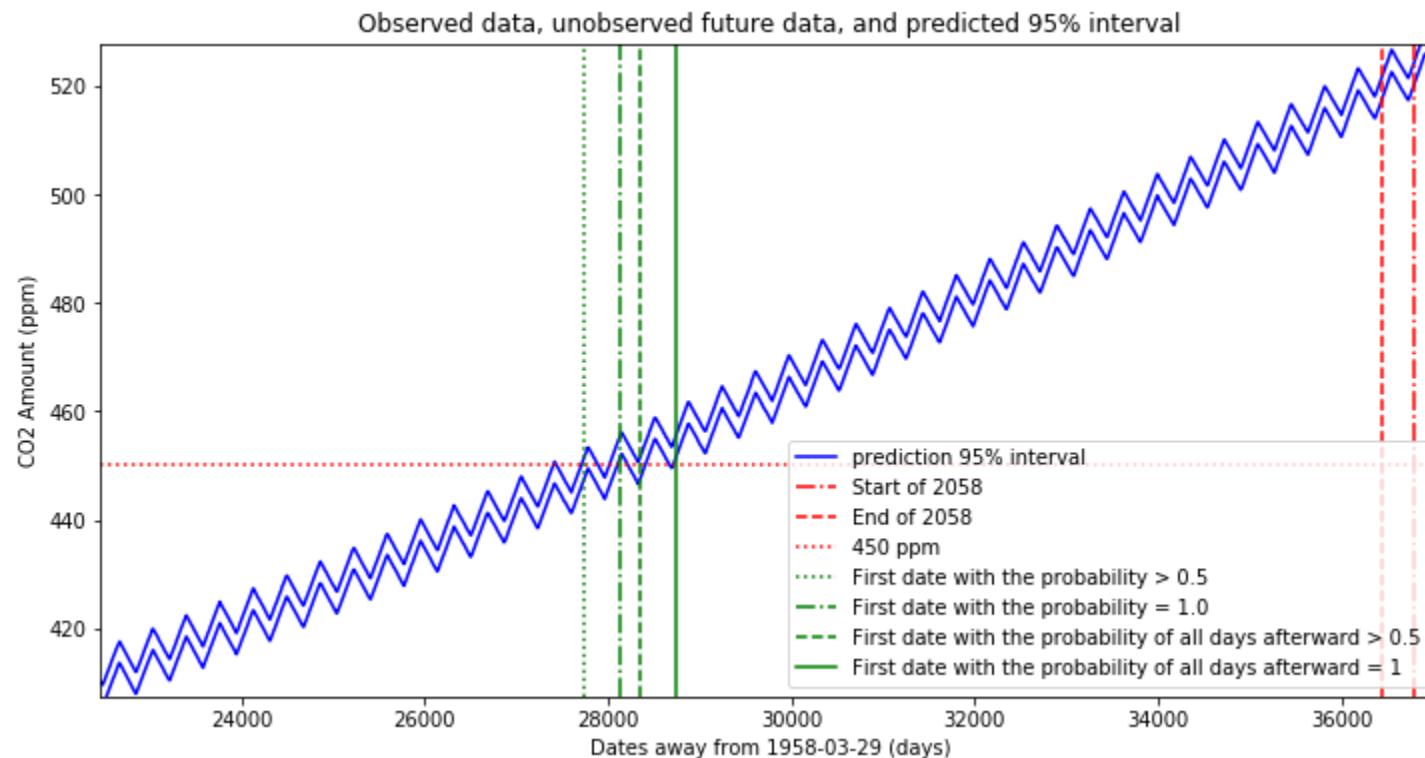


▼ All plot at once

```
1 prediction = samples['x_future']
2
3 # Compute 95% interval of the predicted values
4 prediction_interval = np.percentile(prediction, axis=0, q=[2.5, 97.5])
5
6 # Plot mean and 95% interval of predictions
7 plt.figure(figsize=(12, 6))
8 # plt.plot(data["Time"], data["Amount"], 'k-', label='observed')
9 plt.xlim(data["Time"][-1], max(new_time))
10 plt.ylim(
11     min(prediction[0,:]),
12     max(prediction[1,:]))
13
```

```
14 plt.plot(
15     new_time,
16     prediction_interval[0,:],
17     'b-', label='prediction 95% interval')
18
19 plt.plot(
20     new_time,
21     prediction_interval[1,:],
22     'b-')
23
24 plt.axvline(stop2 * 7, color = "red", linestyle = "-.", label ="Start of 2058")
25 plt.axvline(stop1 * 7, color = "red", linestyle = "--", label ="End of 2058")
26 plt.axhline(450, color = "red", linestyle = ":" , label = "450 ppm")
27 plt.axvline(idx_50 * 7 + data["Time"][n - 1], color = "green", linestyle = ":" , label ="First date with the probability > 0.5")
28 plt.axvline(idx_100 * 7 + data["Time"][n - 1], color = "green", linestyle = "-.", label ="First date with the probability = 1.0")
29 plt.axvline(idx_full_50 * 7 + data["Time"][n - 1], color = "green", linestyle = "--", label ="First date with the probability of al
30 plt.axvline(idx_full_100 * 7 + data["Time"][n - 1], color = "green", linestyle = "-", label ="First date with the probability of al
31
32 plt.xlabel("Dates away from 1958-03-29 (days)")
33 plt.ylabel("CO2 Amount (ppm)")
34 plt.legend()
35 plt.title('Observed data, unobserved future data, and predicted 95% interval')
36 plt.show()
```





1