
An automatic report for the dataset : 02-solar

The Automatic Statistician

Abstract

This report was produced by the Automatic Bayesian Covariance Discovery (ABCD) algorithm.

1 Executive summary

The raw data and full model posterior with extrapolations are shown in figure 1.

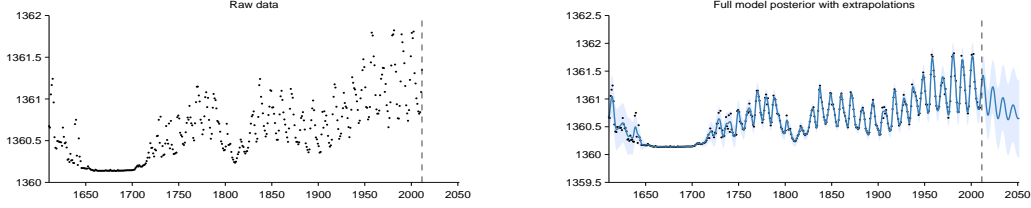


Figure 1: Raw data (left) and model posterior with extrapolation (right)

The structure search algorithm has identified eight additive components in the data. The first 4 additive components explain 92.3% of the variation in the data as shown by the coefficient of determination (R^2) values in table 1. The first 6 additive components explain 99.7% of the variation in the data. After the first 5 components the cross validated mean absolute error (MAE) does not decrease by more than 0.1%. This suggests that subsequent terms are modelling very short term trends, uncorrelated noise or are artefacts of the model or search procedure. Short summaries of the additive components are as follows:

- A constant.
- A constant. This function applies from 1643 until 1716.
- A smooth function. This function applies until 1643 and from 1716 onwards.
- An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards.
- A rapidly varying smooth function. This function applies until 1643 and from 1716 onwards.
- Uncorrelated noise with standard deviation increasing linearly away from 1837. This function applies until 1643 and from 1716 onwards.
- Uncorrelated noise with standard deviation increasing linearly away from 1952. This function applies until 1643 and from 1716 onwards.
- Uncorrelated noise. This function applies from 1643 until 1716.

#	R^2 (%)	ΔR^2 (%)	Residual R^2 (%)	Cross validated MAE	Reduction in MAE (%)
-	-	-	-	1360.65	-
1	0.0	0.0	0.0	0.33	100.0
2	37.4	37.4	37.4	0.23	32.0
3	72.8	35.4	56.6	0.18	21.1
4	92.3	19.4	71.5	0.15	16.8
5	98.1	5.9	75.9	0.15	0.4
6	99.7	1.6	85.6	0.15	0.0
7	100.0	0.3	99.8	0.15	0.0
8	100.0	0.0	100.0	0.15	0.0

Table 1: Summary statistics for cumulative additive fits to the data. The residual coefficient of determination (R^2) values are computed using the residuals from the previous fit as the target values; this measures how much of the residual variance is explained by each new component. The mean absolute error (MAE) is calculated using 10 fold cross validation with a contiguous block design; this measures the ability of the model to interpolate and extrapolate over moderate distances. The model is fit using the full data and the MAE values are calculated using this model; this double use of data means that the MAE values cannot be used reliably as an estimate of out-of-sample predictive performance.

Model checking statistics are summarised in table 2 in section 4. These statistics have revealed statistically significant discrepancies between the data and model in component 8.

The rest of the document is structured as follows. In section 2 the forms of the additive components are described and their posterior distributions are displayed. In section 3 the modelling assumptions of each component are discussed with reference to how this affects the extrapolations made by the model. Section 4 discusses model checking statistics, with plots showing the form of any detected discrepancies between the model and observed data.

2 Detailed discussion of additive components

2.1 Component 1 : A constant

This component is constant.

This component explains 0.0% of the total variance. The addition of this component reduces the cross validated MAE by 100.0% from 1360.6 to 0.3.

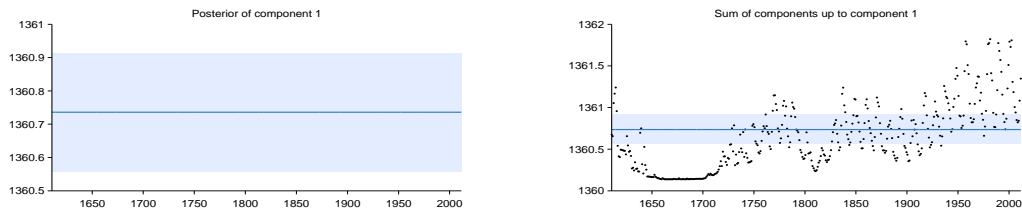


Figure 2: Pointwise posterior of component 1 (left) and the posterior of the cumulative sum of components with data (right)

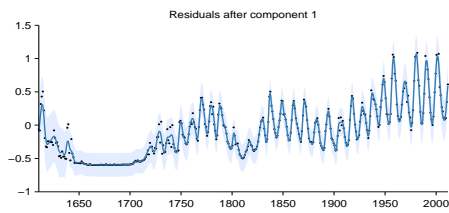


Figure 3: Pointwise posterior of residuals after adding component 1

2.2 Component 2 : A constant. This function applies from 1643 until 1716

This component is constant. This component applies from 1643 until 1716.

This component explains 37.4% of the residual variance; this increases the total variance explained from 0.0% to 37.4%. The addition of this component reduces the cross validated MAE by 31.97% from 0.33 to 0.23.

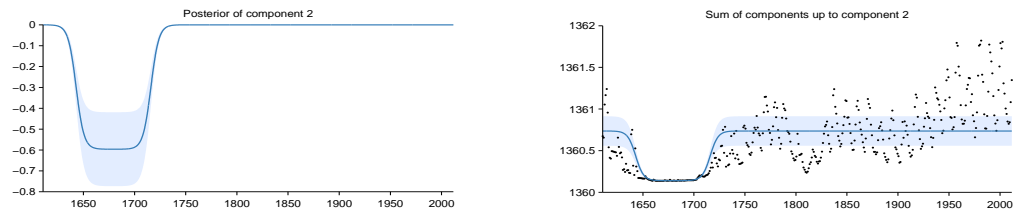


Figure 4: Pointwise posterior of component 2 (left) and the posterior of the cumulative sum of components with data (right)

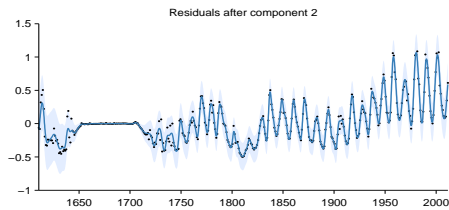


Figure 5: Pointwise posterior of residuals after adding component 2

2.3 Component 3 : A smooth function. This function applies until 1643 and from 1716 onwards

This component is a smooth function with a typical lengthscale of 23.1 years. This component applies until 1643 and from 1716 onwards.

This component explains 56.6% of the residual variance; this increases the total variance explained from 37.4% to 72.8%. The addition of this component reduces the cross validated MAE by 21.08% from 0.23 to 0.18.

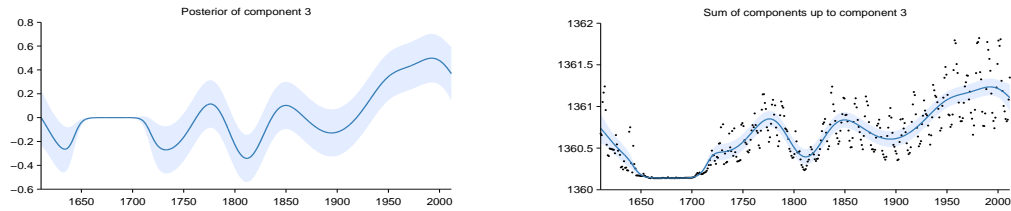


Figure 6: Pointwise posterior of component 3 (left) and the posterior of the cumulative sum of components with data (right)

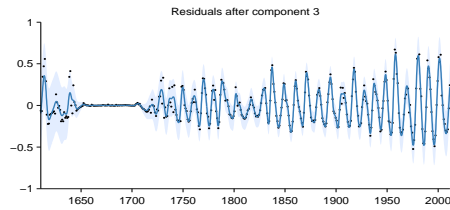


Figure 7: Pointwise posterior of residuals after adding component 3

2.4 Component 4 : An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards

This component is approximately periodic with a period of 10.8 years. Across periods the shape of this function varies smoothly with a typical lengthscale of 36.9 years. The shape of this function within each period is very smooth and resembles a sinusoid. This component applies until 1643 and from 1716 onwards.

This component explains 71.5% of the residual variance; this increases the total variance explained from 72.8% to 92.3%. The addition of this component reduces the cross validated MAE by 16.82% from 0.18 to 0.15.

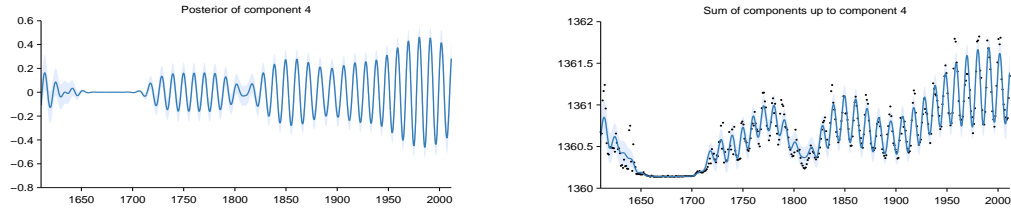


Figure 8: Pointwise posterior of component 4 (left) and the posterior of the cumulative sum of components with data (right)

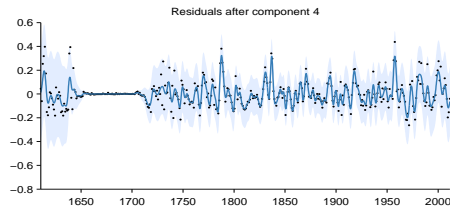


Figure 9: Pointwise posterior of residuals after adding component 4

2.5 Component 5 : A rapidly varying smooth function. This function applies until 1643 and from 1716 onwards

This component is a rapidly varying but smooth function with a typical lengthscale of 1.6 years. This component applies until 1643 and from 1716 onwards.

This component explains 75.9% of the residual variance; this increases the total variance explained from 92.3% to 98.1%. The addition of this component reduces the cross validated MAE by 0.35% from 0.15 to 0.15.

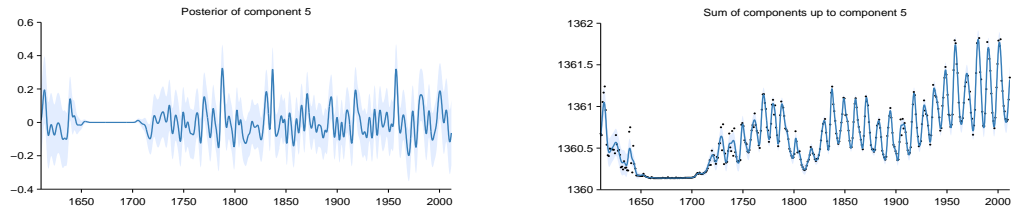


Figure 10: Pointwise posterior of component 5 (left) and the posterior of the cumulative sum of components with data (right)

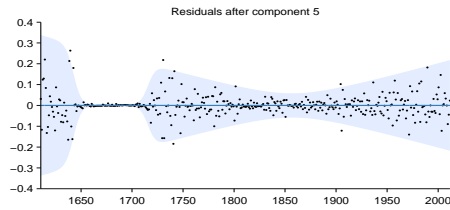


Figure 11: Pointwise posterior of residuals after adding component 5

2.6 Component 6 : Uncorrelated noise with standard deviation increasing linearly away from 1837. This function applies until 1643 and from 1716 onwards

This component models uncorrelated noise. The standard deviation of the noise increases linearly away from 1837. This component applies until 1643 and from 1716 onwards.

This component explains 85.6% of the residual variance; this increases the total variance explained from 98.1% to 99.7%. The addition of this component reduces the cross validated MAE by 0.00% from 0.15 to 0.15. This component explains residual variance but does not improve MAE which suggests that this component describes very short term patterns, uncorrelated noise or is an artefact of the model or search procedure.

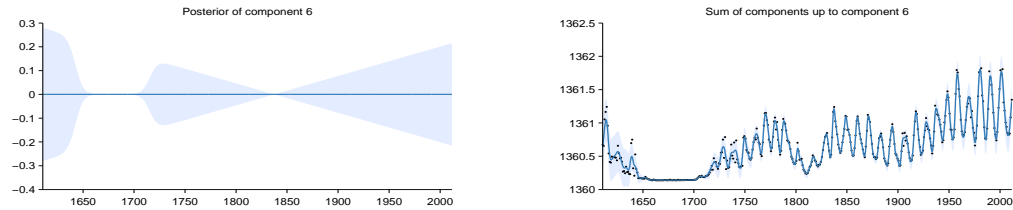


Figure 12: Pointwise posterior of component 6 (left) and the posterior of the cumulative sum of components with data (right)

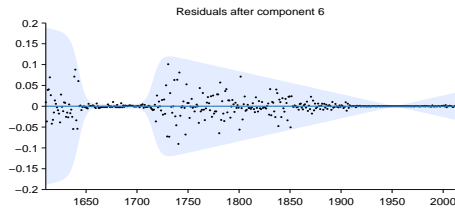


Figure 13: Pointwise posterior of residuals after adding component 6

2.7 Component 7 : Uncorrelated noise with standard deviation increasing linearly away from 1952. This function applies until 1643 and from 1716 onwards

This component models uncorrelated noise. The standard deviation of the noise increases linearly away from 1952. This component applies until 1643 and from 1716 onwards.

This component explains 99.8% of the residual variance; this increases the total variance explained from 99.7% to 100.0%. The addition of this component reduces the cross validated MAE by 0.00% from 0.15 to 0.15. This component explains residual variance but does not improve MAE which suggests that this component describes very short term patterns, uncorrelated noise or is an artefact of the model or search procedure.

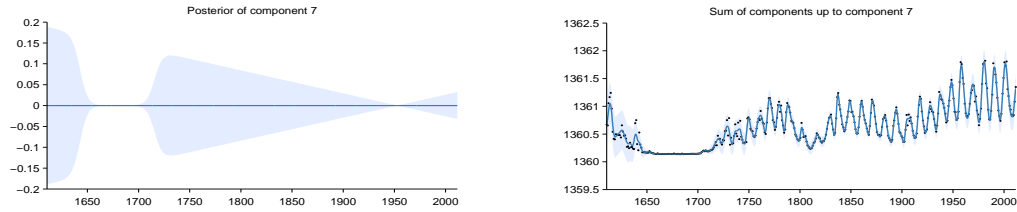


Figure 14: Pointwise posterior of component 7 (left) and the posterior of the cumulative sum of components with data (right)

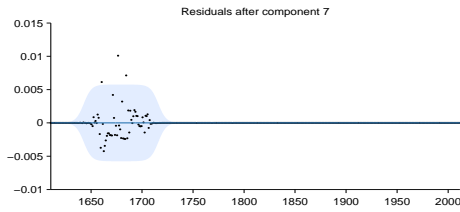


Figure 15: Pointwise posterior of residuals after adding component 7

2.8 Component 8 : Uncorrelated noise. This function applies from 1643 until 1716

This component models uncorrelated noise. This component applies from 1643 until 1716.

This component explains 100.0% of the residual variance; this increases the total variance explained from 100.0% to 100.0%. The addition of this component reduces the cross validated MAE by 0.00% from 0.15 to 0.15. This component explains residual variance but does not improve MAE which suggests that this component describes very short term patterns, uncorrelated noise or is an artefact of the model or search procedure.

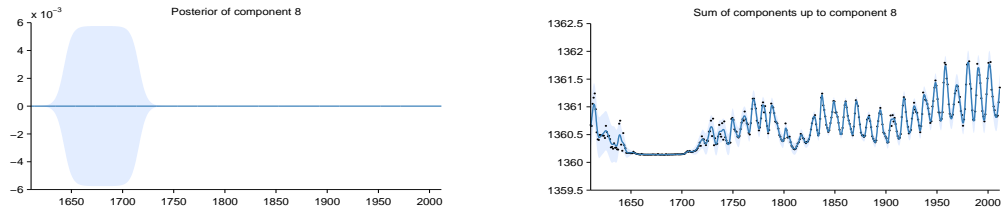


Figure 16: Pointwise posterior of component 8 (left) and the posterior of the cumulative sum of components with data (right)

3 Extrapolation

Summaries of the posterior distribution of the full model are shown in figure 17. The plot on the left displays the mean of the posterior together with pointwise variance. The plot on the right displays three random samples from the posterior.

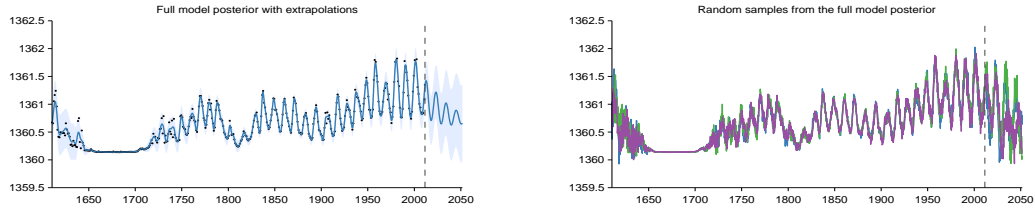


Figure 17: Full model posterior with extrapolation. Mean and pointwise variance (left) and three random samples (right)

Below are descriptions of the modelling assumptions associated with each additive component and how they affect the predictive posterior. Plots of the pointwise posterior and samples from the posterior are also presented, showing extrapolations from each component and the cumulative sum of components.

3.1 Component 1 : A constant

This component is assumed to stay constant.

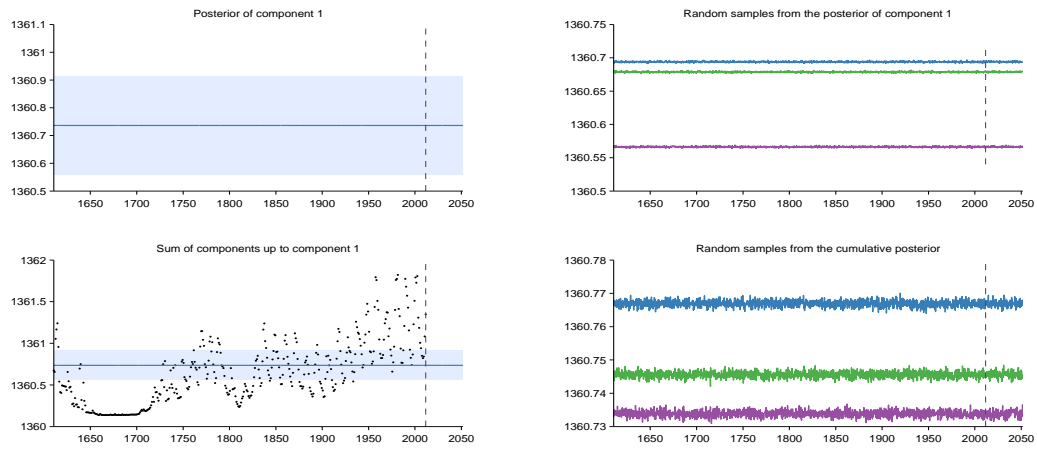


Figure 18: Posterior of component 1 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.2 Component 2 : A constant. This function applies from 1643 until 1716

This component is assumed to stop before the end of the data and will therefore be extrapolated as zero.

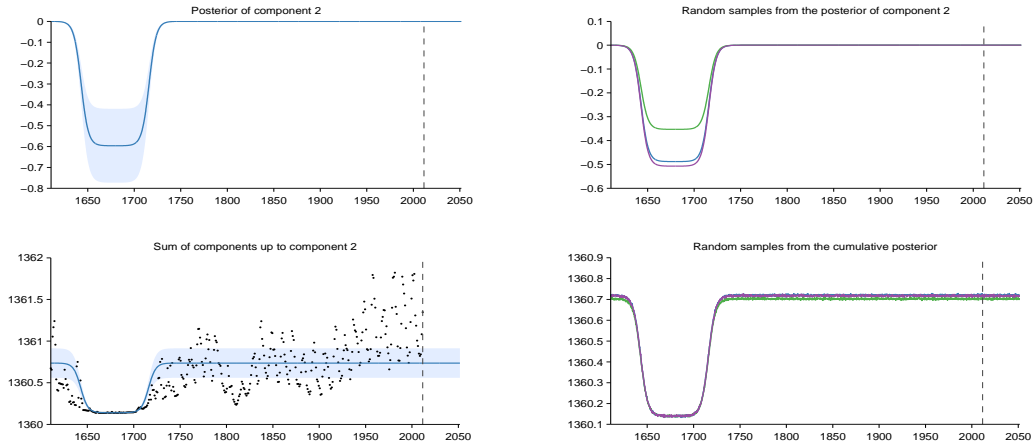


Figure 19: Posterior of component 2 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.3 Component 3 : A smooth function. This function applies until 1643 and from 1716 onwards

This component is assumed to continue smoothly but is also assumed to be stationary so its distribution will return to the prior. The prior distribution places mass on smooth functions with a marginal mean of zero and a typical lengthscale of 23.1 years. [This is a placeholder for a description of how quickly the posterior will start to resemble the prior].

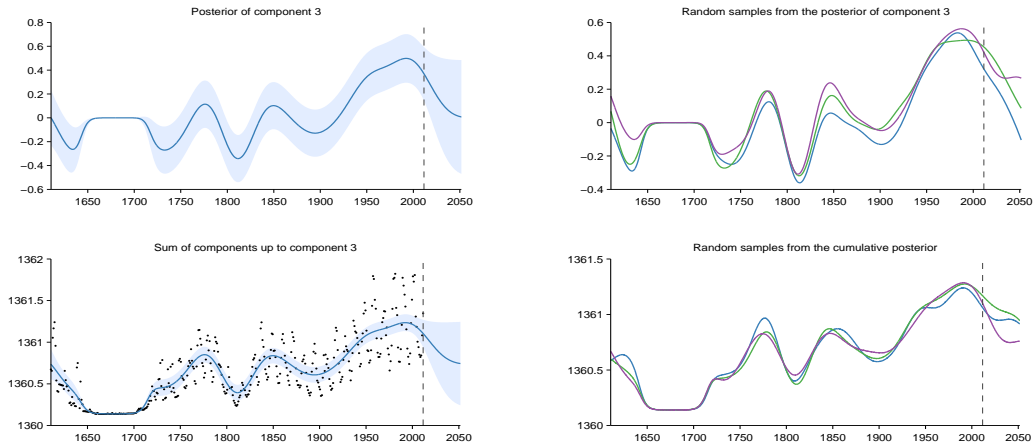


Figure 20: Posterior of component 3 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.4 Component 4 : An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards

This component is assumed to continue to be approximately periodic. The shape of the function is assumed to vary smoothly between periods but will return to the prior. The prior is entirely uncertain about the phase of the periodic function. Consequently the pointwise posterior will appear to lose its periodicity, but this merely reflects the uncertainty in the shape and phase of the function. [This is a placeholder for a description of how quickly the posterior will start to resemble the prior].

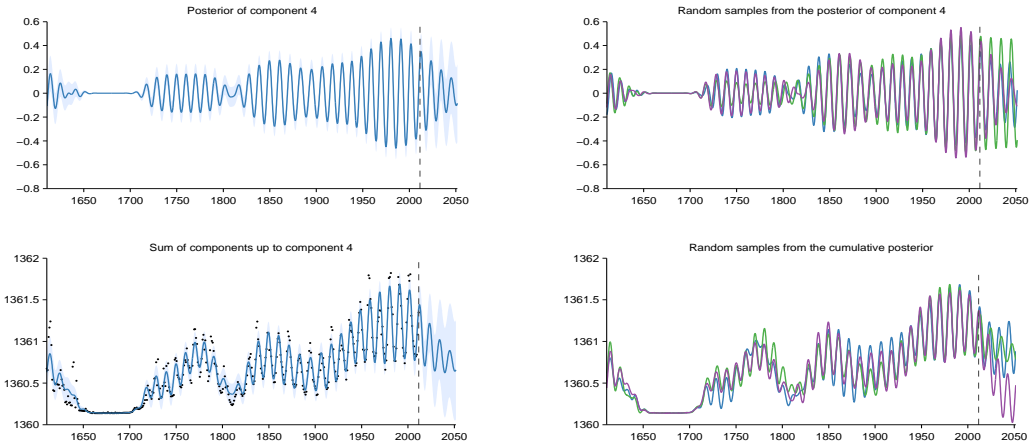


Figure 21: Posterior of component 4 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.5 Component 5 : A rapidly varying smooth function. This function applies until 1643 and from 1716 onwards

This component is assumed to continue smoothly but its distribution is assumed to quickly return to the prior. The prior distribution places mass on smooth functions with a marginal mean of zero and a typical lengthscale of 1.6 years. [This is a placeholder for a description of how quickly the posterior will start to resemble the prior].

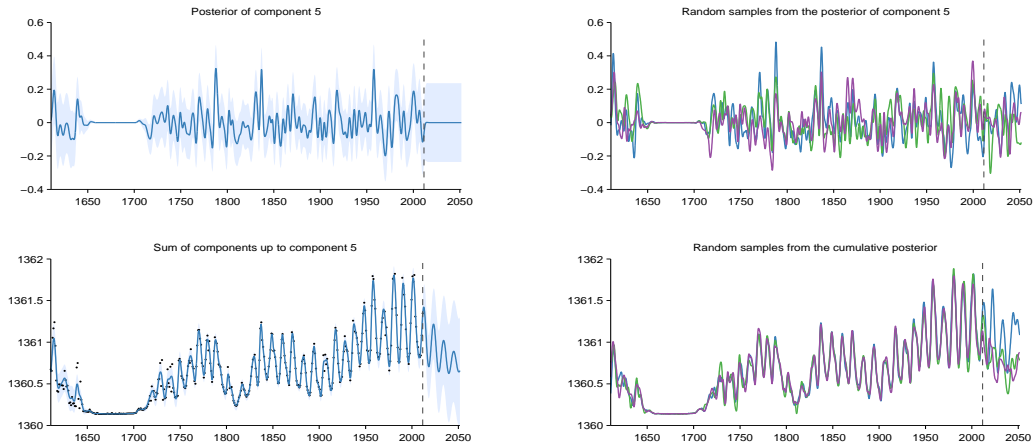


Figure 22: Posterior of component 5 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.6 Component 6 : Uncorrelated noise with standard deviation increasing linearly away from 1837. This function applies until 1643 and from 1716 onwards

This component assumes the uncorrelated noise will continue indefinitely. The standard deviation of the noise is assumed to continue to increase linearly.

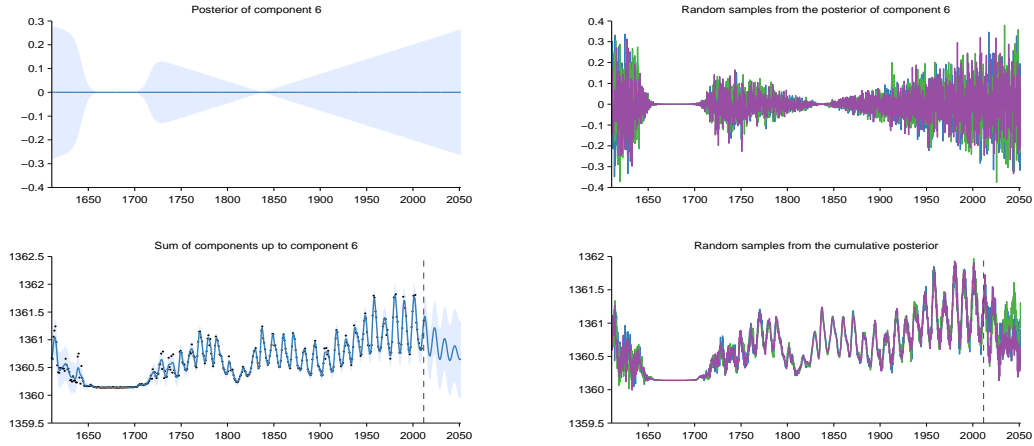


Figure 23: Posterior of component 6 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.7 Component 7 : Uncorrelated noise with standard deviation increasing linearly away from 1952. This function applies until 1643 and from 1716 onwards

This component assumes the uncorrelated noise will continue indefinitely. The standard deviation of the noise is assumed to continue to increase linearly.

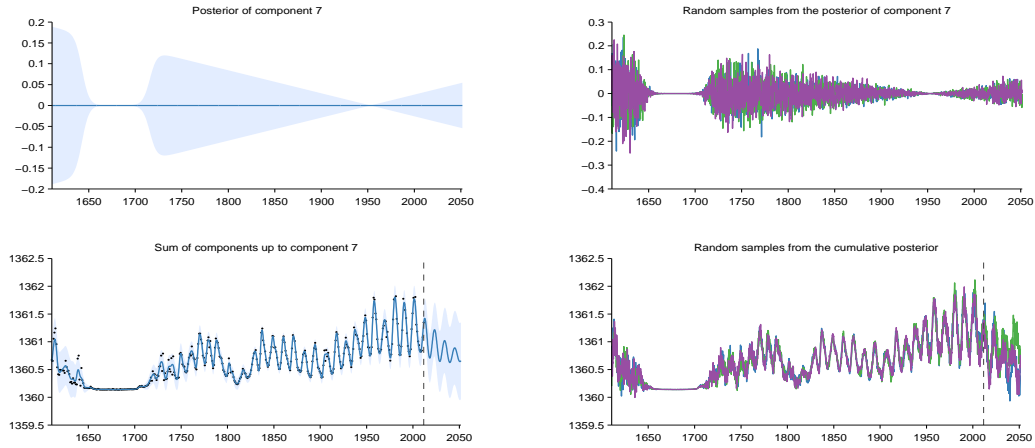


Figure 24: Posterior of component 7 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

3.8 Component 8 : Uncorrelated noise. This function applies from 1643 until 1716

This component is assumed to stop before the end of the data and will therefore be extrapolated as zero.

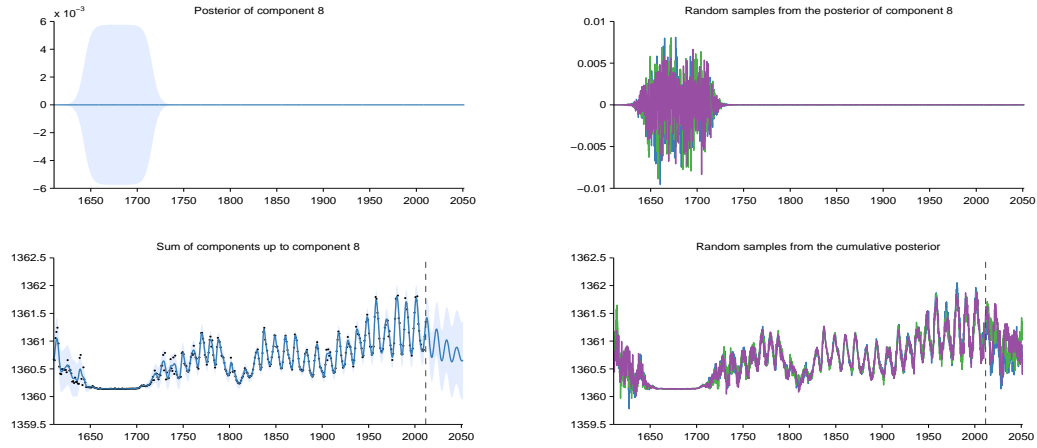


Figure 25: Posterior of component 8 (top) and cumulative sum of components (bottom) with extrapolation. Mean and pointwise variance (left) and three random samples from the posterior distribution (right).

4 Model checking

Several posterior predictive checks have been performed to assess how well the model describes the observed data. These tests take the form of comparing statistics evaluated on samples from the prior and posterior distributions for each additive component. The statistics are derived from autocorrelation function (ACF) estimates, periodograms and quantile-quantile (qq) plots.

Table 2 displays cumulative probability and p -value estimates for these quantities. Cumulative probabilities near 0/1 indicate that the test statistic was lower/higher under the posterior compared to the prior unexpectedly often i.e. they contain the same information as a p -value for a two-tailed test and they also express if the test statistic was higher or lower than expected. p -values near 0 indicate that the test statistic was larger in magnitude under the posterior compared to the prior unexpectedly often.

#	ACF		Periodogram		QQ	
	min	min loc	max	max loc	max	min
1	0.501	0.481	0.543	0.497	0.223	0.776
2	0.501	0.479	0.723	0.500	0.858	0.192
3	0.959	0.898	0.734	0.229	0.368	0.792
4	0.564	0.486	0.393	0.371	0.790	0.812
5	0.605	0.465	0.409	0.455	0.204	0.732
6	0.516	0.477	0.412	0.396	0.477	0.674
7	0.456	0.510	0.461	0.480	0.498	0.561
8	0.584	0.638	0.585	0.526	0.012	0.697

Table 2: Model checking statistics for each component. Cumulative probabilities for minimum of autocorrelation function (ACF) and its location. Cumulative probabilities for maximum of periodogram and its location. p -values for maximum and minimum deviations of QQ-plot from straight line.

The nature of any observed discrepancies is now described and plotted and hypotheses are given for the patterns in the data that may not be captured by the model.

4.1 Moderately statistically significant discrepancies

4.1.1 Component 8 : Uncorrelated noise. This function applies from 1643 until 1716

The following discrepancies between the prior and posterior distributions for this component have been detected.

- The qq plot has an unexpectedly large positive deviation from equality ($x = y$). This discrepancy has an estimated p -value of 0.012.

The positive deviation in the qq-plot can indicate heavy positive tails if it occurs at the right of the plot or light negative tails if it occurs as the left.

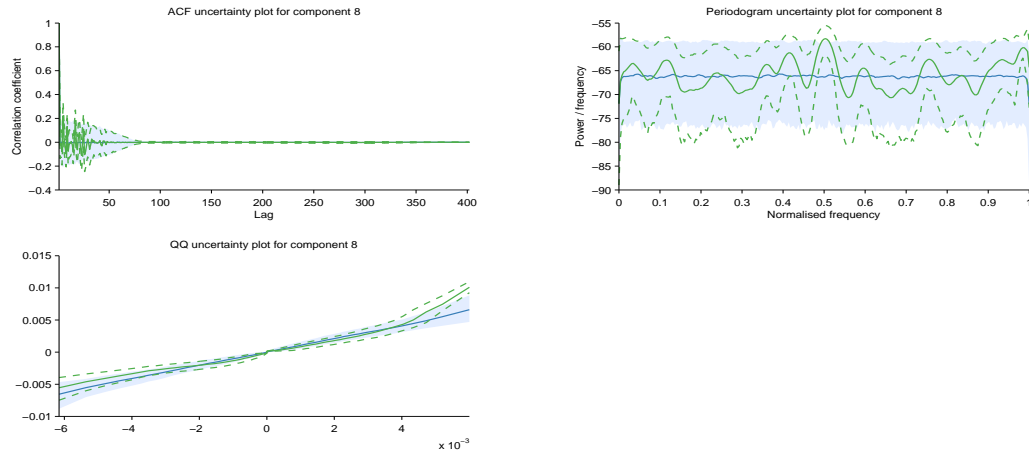


Figure 26: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 8. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2 Model checking plots for components without statistically significant discrepancies

4.2.1 Component 1 : A constant

No discrepancies between the prior and posterior of this component have been detected

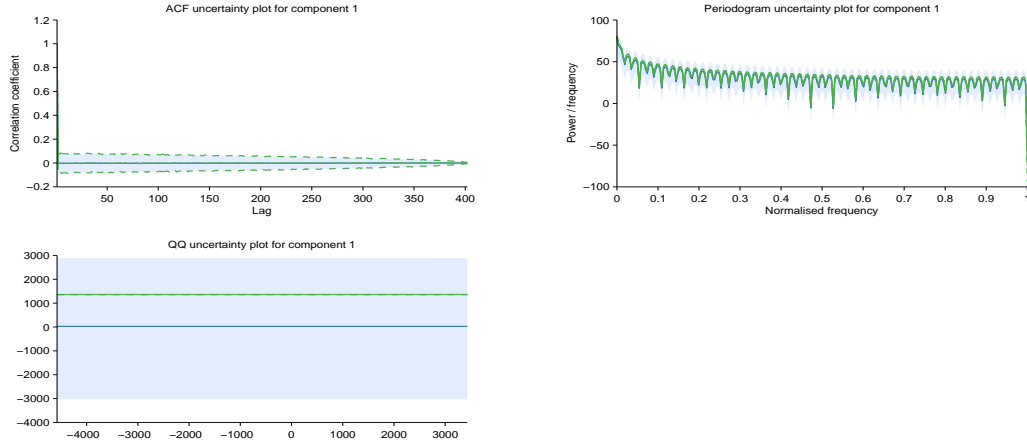


Figure 27: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 1. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.2 Component 2 : A constant. This function applies from 1643 until 1716

No discrepancies between the prior and posterior of this component have been detected

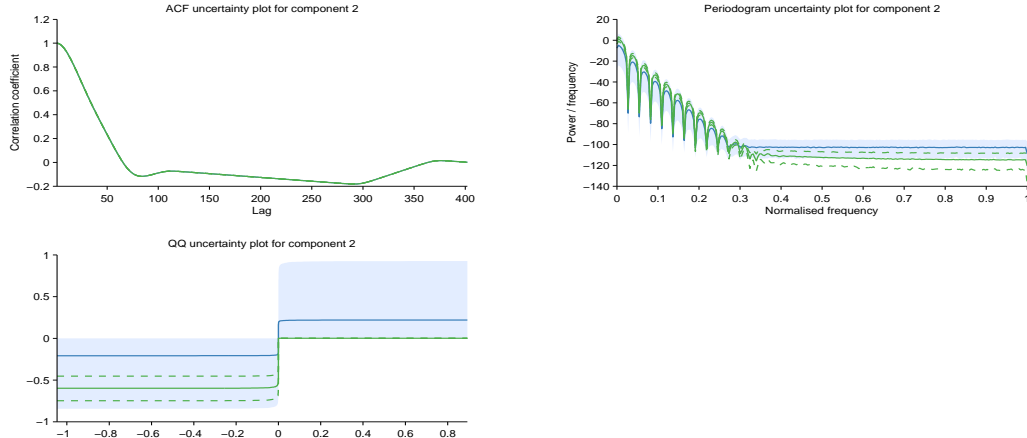


Figure 28: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 2. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.3 Component 3 : A smooth function. This function applies until 1643 and from 1716 onwards

No discrepancies between the prior and posterior of this component have been detected

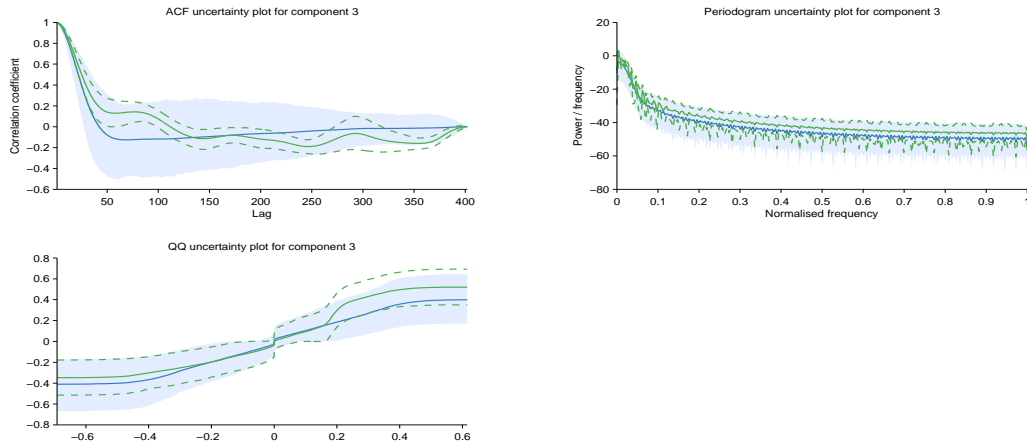


Figure 29: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 3. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.4 Component 4 : An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards

No discrepancies between the prior and posterior of this component have been detected

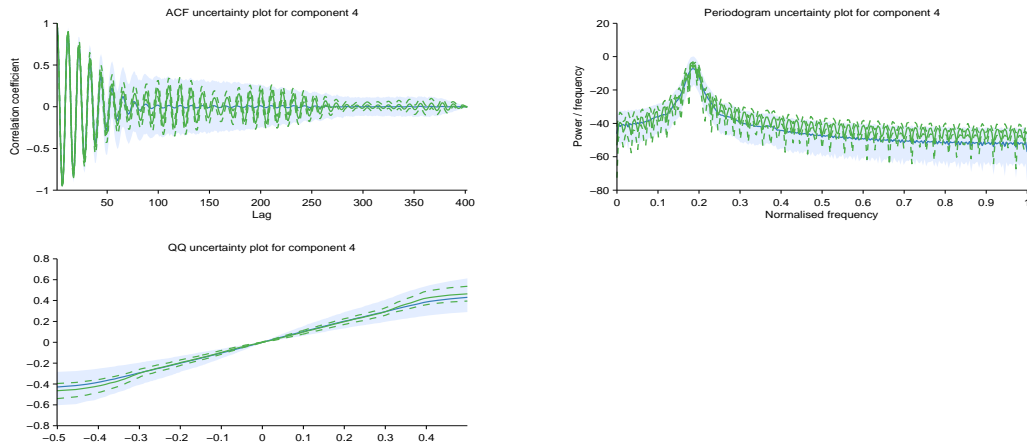


Figure 30: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 4. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.5 Component 5 : A rapidly varying smooth function. This function applies until 1643 and from 1716 onwards

No discrepancies between the prior and posterior of this component have been detected

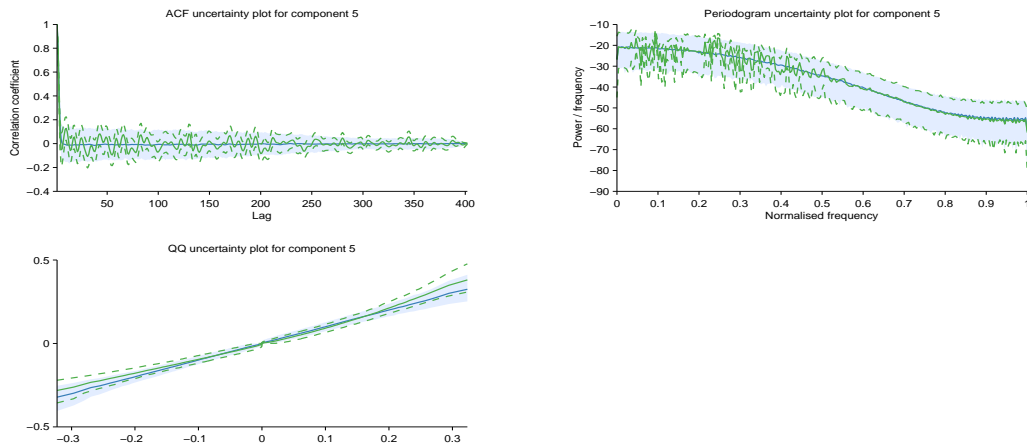


Figure 31: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 5. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.6 Component 6 : Uncorrelated noise with standard deviation increasing linearly away from 1837. This function applies until 1643 and from 1716 onwards

No discrepancies between the prior and posterior of this component have been detected

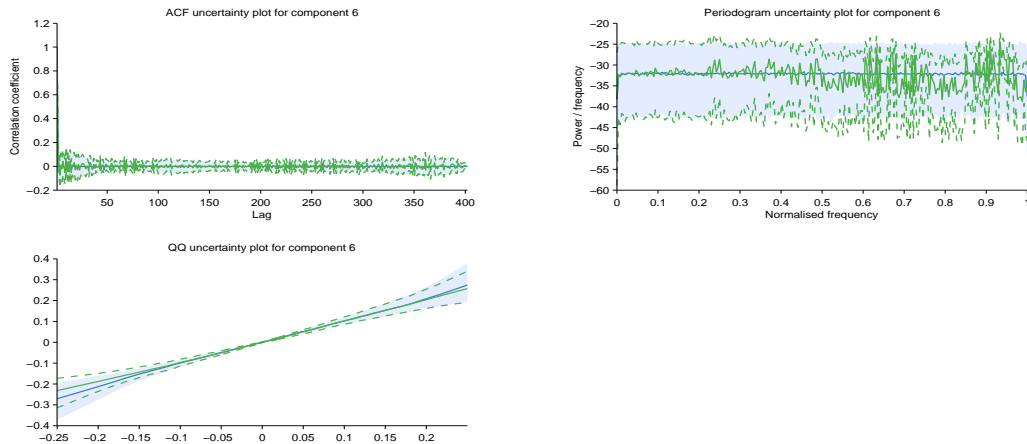


Figure 32: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 6. The green line and green dashed lines are the corresponding quantities under the posterior.

4.2.7 Component 7 : Uncorrelated noise with standard deviation increasing linearly away from 1952. This function applies until 1643 and from 1716 onwards

No discrepancies between the prior and posterior of this component have been detected

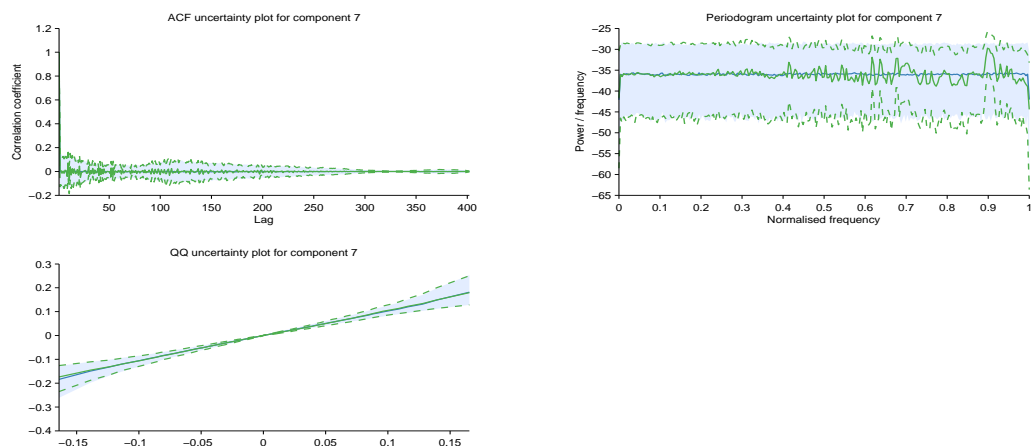


Figure 33: ACF (top left), periodogram (top right) and quantile-quantile (bottom left) uncertainty plots. The blue line and shading are the pointwise mean and 90% confidence interval of the plots under the prior distribution for component 7. The green line and green dashed lines are the corresponding quantities under the posterior.