

Global CO_2 Emissions in 1997

Naikaj Pandya, Amit Gattadahalli, Michael Golas, Austin Pitts

Introduction

In 1960, Charles Keeling in his seminal paper, *The Concentration and Isotopic Abundances of Carbon Dioxide in the Atmosphere* made two notable observations:

1. That a seasonal variation in CO_2 concentrations was observed in the northern hemisphere, corresponding to the activity of land plants
2. That at longer horizons, beyond one year, global concentrations of CO_2 have increased at a rate of 1.3 p.p.m. either from the combustion of fossil fuels or from factors tied to the seasonal variation, exceeding the counteracting oceanic effect removing CO_2 from the atmosphere.

Keeling's analysis was conducted using data obtained from three gas analyzers, equipped to measure carbon dioxide concentrations continuously, located in Antarctica, Hawaii and California.

Our goal, in 1994, is to validate Keeling's observations using data collected during the intervening years, measured using modern optical sensors at higher frequencies and report any observed changes to the rates of accumulated CO_2 in the atmosphere. Using these estimates we plan to extend our study and apply time-series modeling techniques to forecast the trends and variation in expected future CO_2 concentrations to provide bounds on the anticipated levels of CO_2 . Since the amount of atmospheric CO_2 carries broad environmental and economic effects, our results are relevant to both environmental and policy researchers as crucial estimates to help guide mitigating courses of action within the appropriate time frames.

CO2 Data

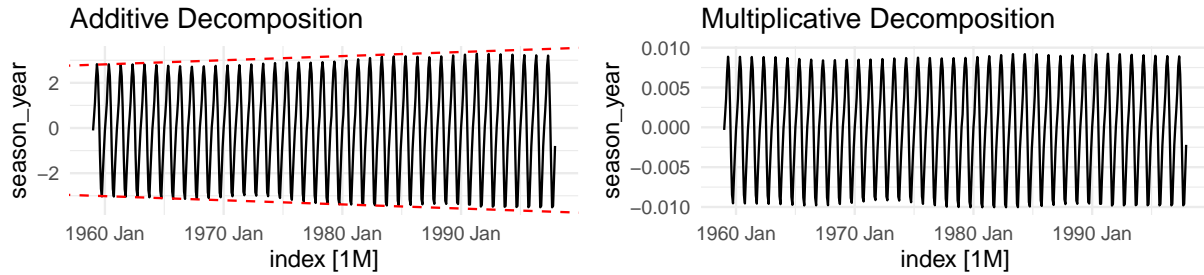
We conduct a timeseries-analysis of atmospheric CO_2 levels using data collected by the NOAA. As stated, the goal is to examine any long-term trends and seasonal fluctuations in CO_2 levels.

The input data sourced from NOAA is collected using a CO_2 analyzer installed at Mauna Loa that uses a technique based on infrared absorption, wherein a sensor measures the magnitude of absorption of light circulating in an optical cavity. Data is collected hourly, daily and monthly, we use the monthly average data for this analysis as our primary interest is devoted to long-term changes in CO_2 levels. An important aspect of the measurements is the ongoing calibrations of the analyzer. The absorption by the instrument depends on the total amount of CO_2 , therefore the temperature and pressure in the instrument, as well as the flow rate, need to be measured and frequent calibrations performed with reference gas mixtures of known amounts of CO_2 -in-dry-air. The intake lines are from the top of a 38 m tall tower next to the observatory, to avoid any influence on the measurements by human activities at the observatory. The difference of the ambient air measurements from the reference gas R0 are calculated, and these differences are used to calculate the true fraction CO_2 .

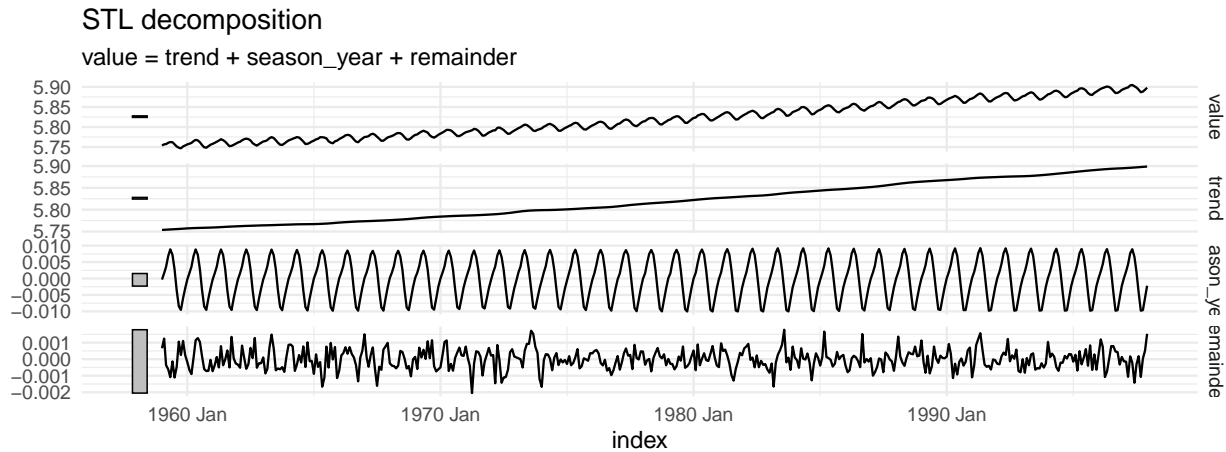
We begin with an exploratory analysis of the data guided by a few general observations apparent from visual inspection of the time-series:

- the data shows variation periodic in time
- the general level increases over time

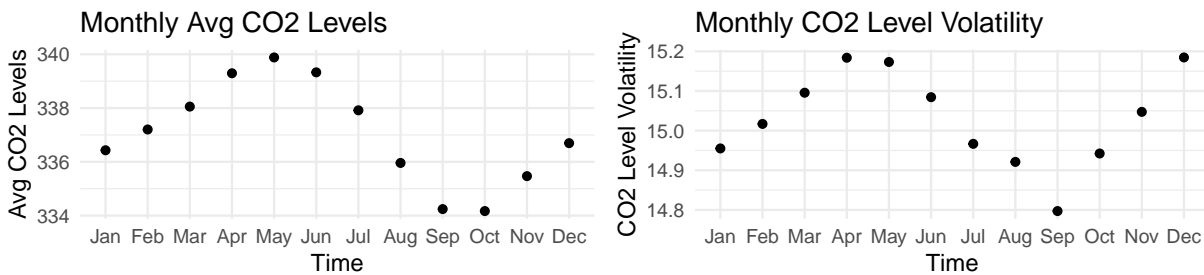
One question is whether the variation remains constant independent of the level of CO_2 . Analysis of variation around the trend-cycle reveals a persistent increase in the amplitude of the fluctuations.



As we can see in the figures above, the Additive Decomposition flares outward. This leads us to conclude that the appropriate decomposition of the time-series into Trend, Seasonal, and Residual components is via Multiplicative Decomposition.



Looking at the STL decomposition, although the long run growth rate of co_2 is very low, approx. 0.0127402% per year, almost linear at the time-scale of observation, we note that the growth is highly statistically significant.



Finally, observing the month-to-month average CO_2 levels and volatility gives us an idea of the seasonality in our data.

Linear Time Trend Model

To setup our problem for validation we split our dataset into an in-sample train period spanning years prior to 1997 and post-1997 as the test period.

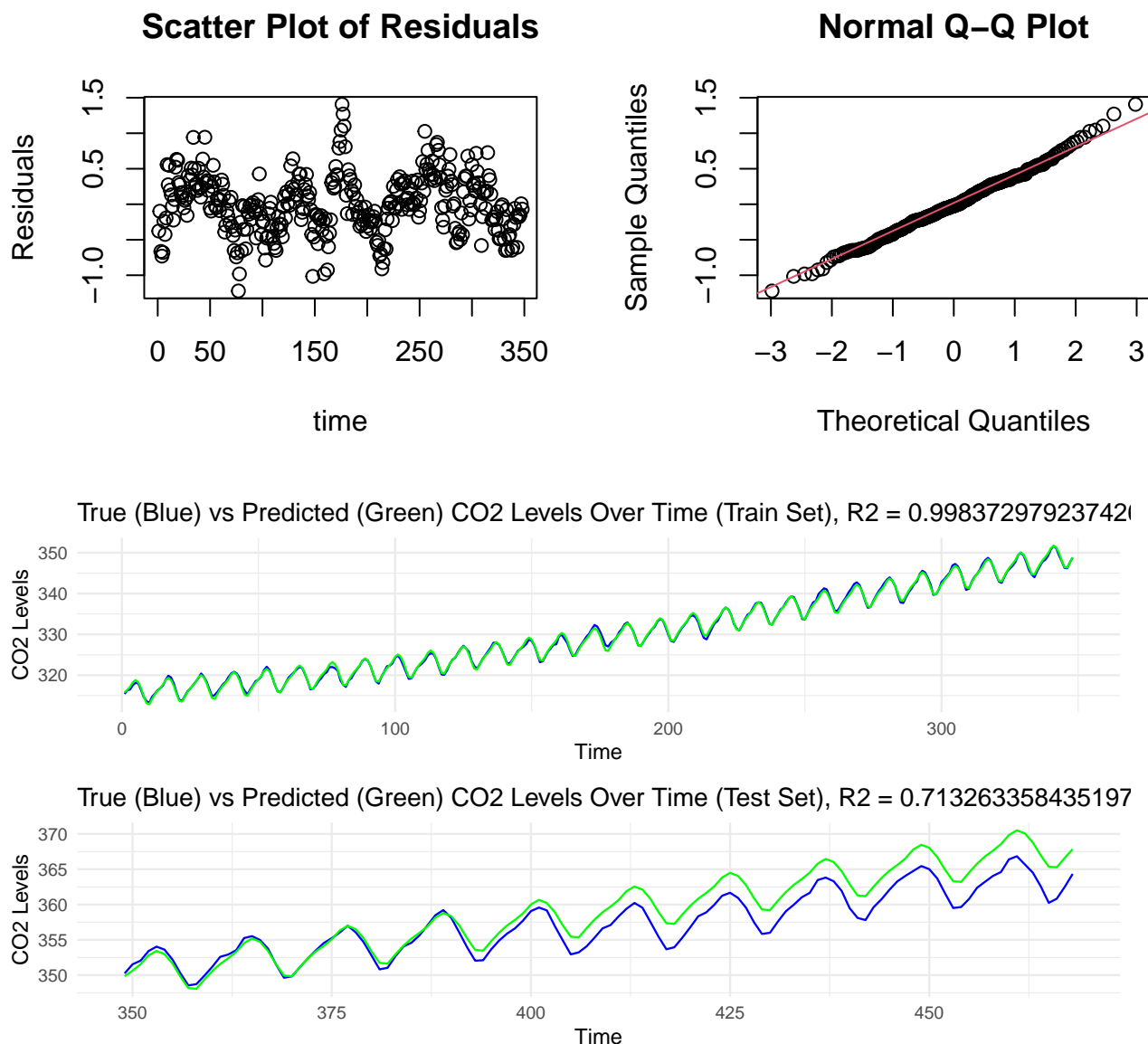
We fit linear, polynomial and quadratic models. We observe that the using a linear timeseries model the residuals exhibit a positive trend. This is also prevalent for the quadratic and polynomial models, however the magnitude of the trend is reduced. Our final model uses linear, quadratic, exponential and seasonal features.

```
final_model = lm(value~index1+index2+log_index + month,df)
```

For each of our models we evaluated model performance by observing the following visualizations:

- Scatter Plot of Residuals
- Normal Q-Q Plot
- True (Blue) vs Predicted (Green) CO2 Levels Over Time (Train Set)
- True (Blue) vs Predicted (Green) CO2 Levels Over Time (Test Set)

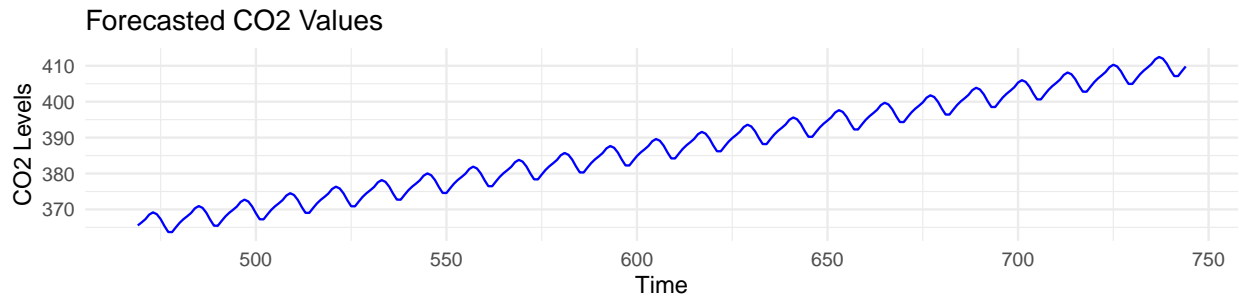
Below are the plots generated by our final model when trained using our train data. For the plots generated by each of the models we analyzed before coming to this final model, see Appendix for the full notebook.



From the above visualizations of our final model we can observe the following. The qq plot mostly resembles a straight line with less deviation than the qq plot of the other models we tried, suggesting that our data is approximately normally distributed. For the train set, our true vs predicted values of CO2 levels over time stayed very close to each other. Finally, for the test set, our true vs predicted values of CO2 levels over

time mimicked the seasonality and growth but slowly our predicted separated from the true values as large amounts of time passed. Overall, we found that this was the model that produced the best results.

We then used our final linear model to generate forecasts to the year 2022.



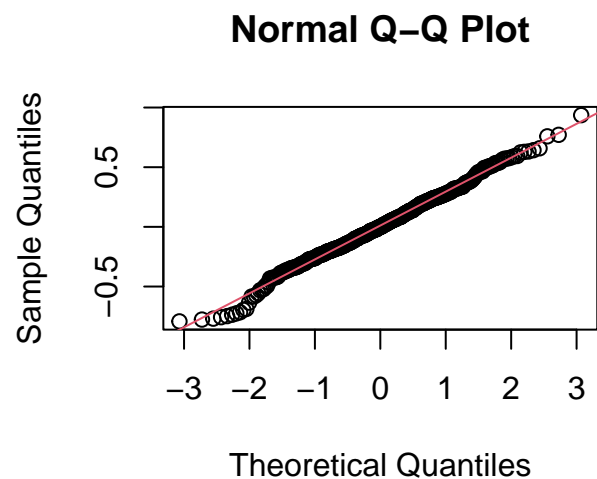
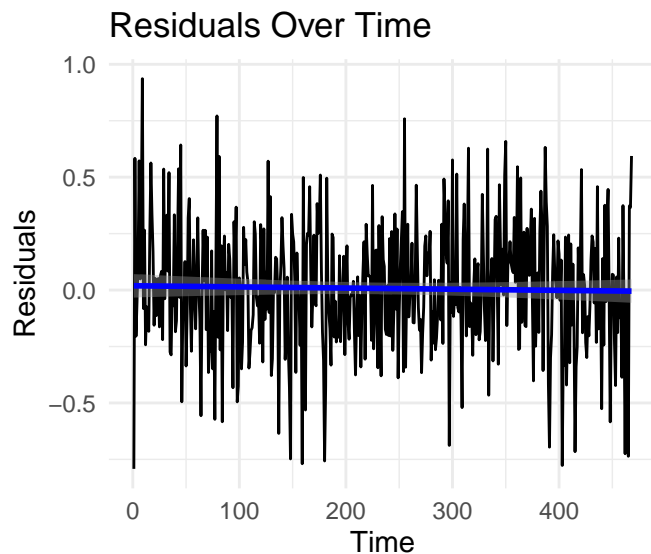
ARIMA Time Series Model

We developed an ARIMA model to fit to the series and generate forecasts to the year 2022. After finding the optimal linear model from the previous section, we trained a linear model on the full dataset and used it to detrend out data prior to arima model training. We find that it is necessary to detrend the CO_2 Series such that it is stationary. To achieve this we use a linear model to detrend series as it captures the linear & nonlinear temporal and seasonal trends inherent in the data.

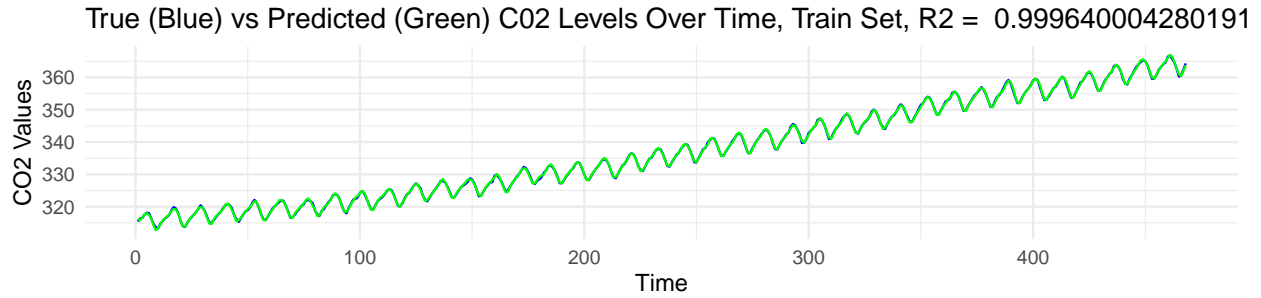
```
df2 = df
df2[['diff_value']] = final_model$residuals
arima_model = arima(df2[['diff_value']], order=c(2,0,0))
```

While there is not enough room to provide all of these plots in this report (see Appendix for full notebook), in doing this process, we also

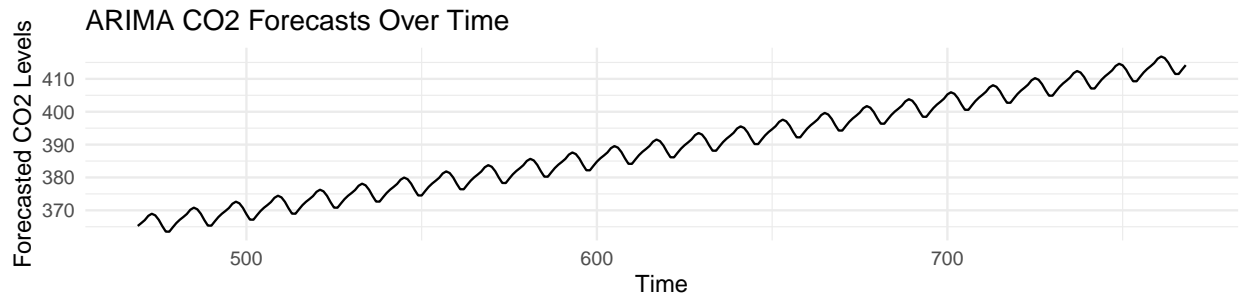
- checked mean/variance of series over time to visually validate stationarity assumptions
- checked ACF/PACF plots
- performed a grid search to find P, Q values that minimize in sample BIC
- validated that our final model's residuals are white noise and approximately normally distributed (shown below)



We also plotted true vs predicted CO_2 levels over time for evaluation of our ARIMA model. As can be seen in the below plot, the predicted very closely followed the true values.

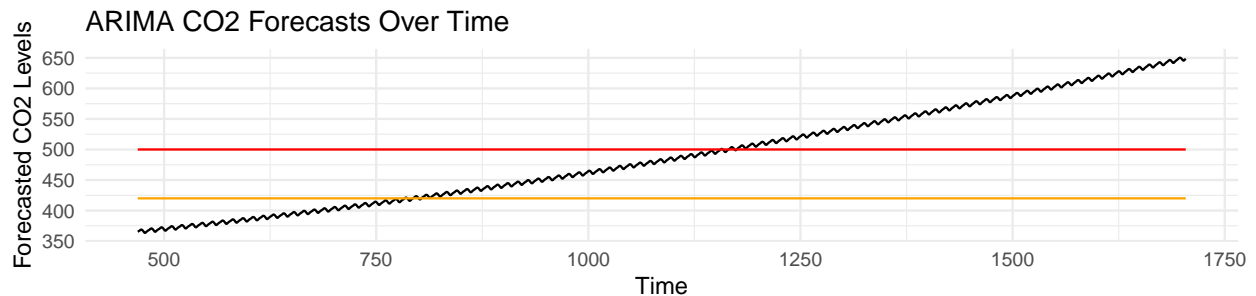


Similar to our linear model, we then used our final ARIMA model to also generate forecasts to the year 2022.



Forecasting Atmospheric CO2 Growth

We generate predictions for when atmospheric CO2 is expected to be at 420ppm and 500 ppm levels for the first and final times. Given errors observed from previous forecasts our hope is that our more modern models which utilize more sophisticated features will provide more accurate results of future CO_2 levels.



Based on this forecasting using our ARIMA model, our predictions are as follows:

- First and Final Time at 420, April 2024 - Oct 2026
- First and Final Time at 500, April 2055 - Nov 2056

We are fairly confident that these will be close to accurate predictions based on our analysis of our ARIMA model and its performance, but forecasting so far into the future means that our predictions will likely not be perfect.

Appendix

While our final results are reported here, in our complete notebook (Github Folder: Notebook) we examine alternative models and go into further assessment of the models. The purpose of this background information is to show more of the process in how we reached the conclusions shown above.

Global CO_2 Emissions in the Present

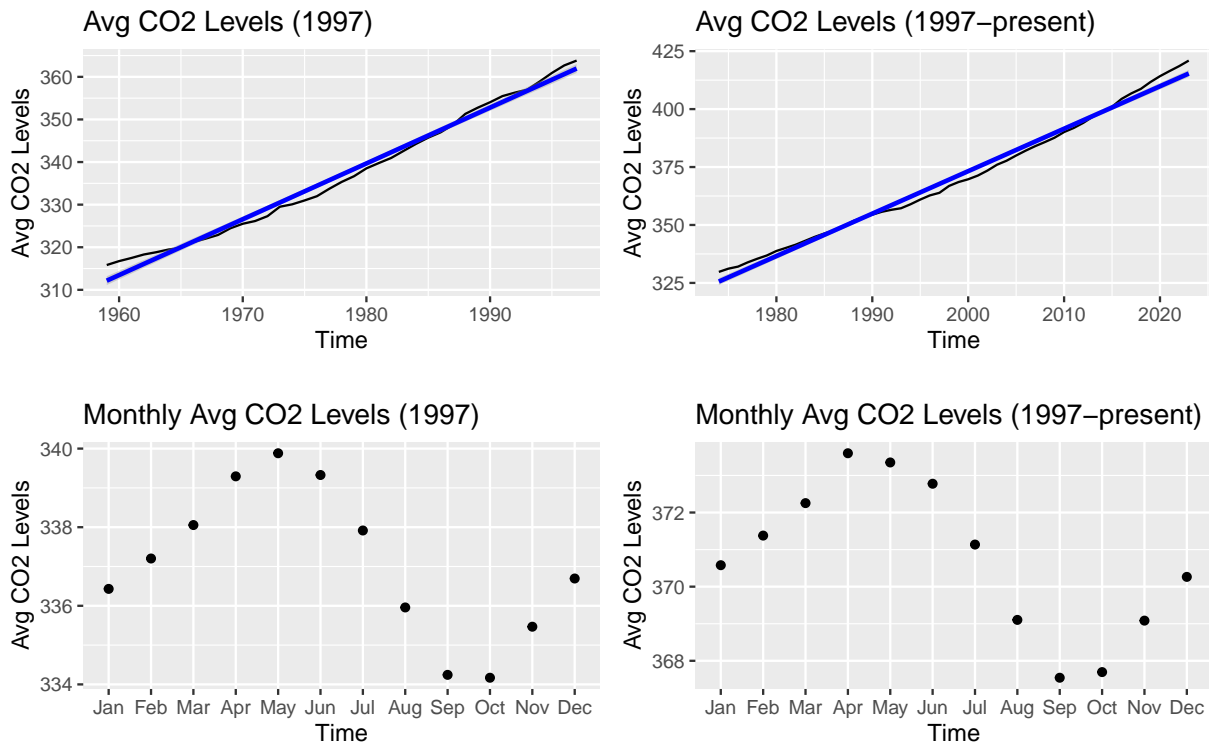
Naikaj Pandya, Amit Gattadahalli, Michael Golas, Austin Pitts

Introduction

Our goal is to re-evaluate Keeling's observations of accumulated CO_2 in the atmosphere using the most upto date data. As of April of 2019, a new CO_2 analyzer was installed at Mauna Loa that uses a technique called Cavity Ring-Down Spectroscopy (CRDS). CRDS is based on the measurement of the rate of absorption of light circulating in an optical cavity by comparing the ring down times when the laser is at a wavelength that the CO_2 molecule does not absorb, to the ring down time when the laser is at a wavelength that the CO_2 molecule does absorb.

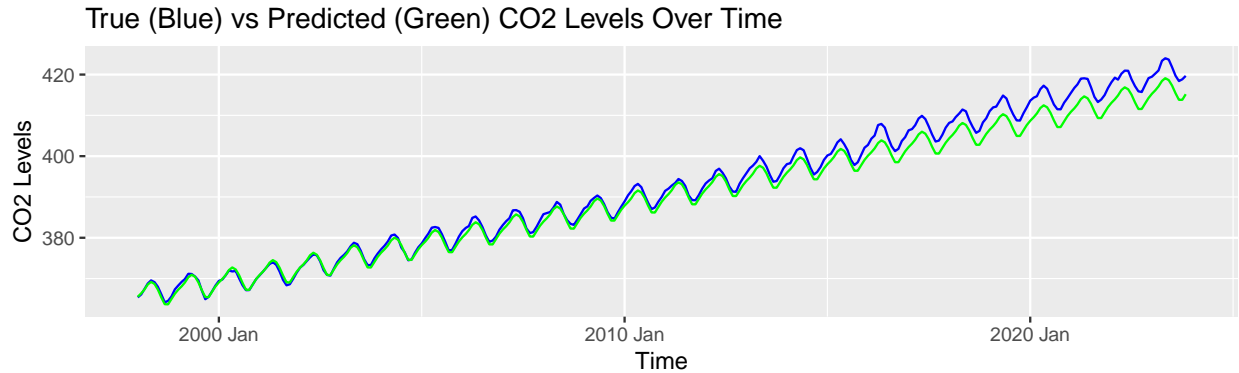
Create a Modern Data Pipeline for Mona Loa CO_2 Data

We establish a modern data pipeline for Mona Loa CO_2 data. The data was obtained from the NOAA website, specifically from the CO_2 daily data page. Additional columns for analysis were created, such as a log-transformed index and polynomial terms for time. The time series data was converted into a tsibble object for efficient time-series analysis. We then conducted an exploratory data analysis (EDA) to gain insights into the CO_2 levels over time. The analysis included visualizations depicting the overall trend, average/standard deviation (SD) of CO_2 levels per year, and monthly variations. The first plot illustrated the continuous increase in CO_2 levels over the years prior to 1997 then the second plot shows from 1997 to present. The third and forth plots delved into monthly variations, highlighting average CO_2 levels and their volatility. This EDA lays the foundation for further analysis and comparison with Keeling's observations from 1997. It also sets the stage for evaluating the performance of earlier models and forecasting future CO_2 levels.



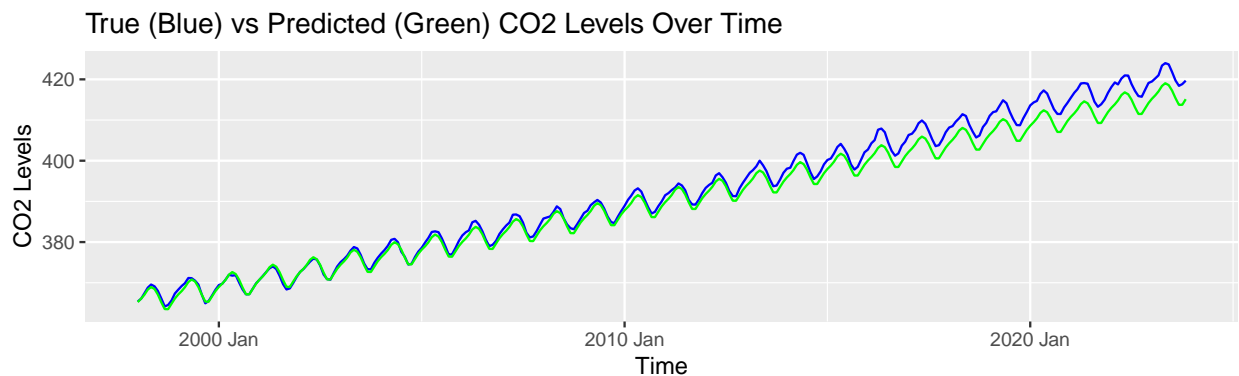
Compare Linear Model Forecasts Against Realized CO2

In our comparison of realized atmospheric CO2 levels to those predicted by your forecast from a linear time model in 1997 (i.e. “Task 2a”), we can see the results in our graph of True (Blue) vs Predicted (Green) CO2 Levels Over Time. Results stay very consistent at first but as time continues we can see that the true values are increasing at a larger rate than what our predicted values would have estimated.



Compare ARIMA Model Forecasts Against Realized CO2

In our comparison of realized atmospheric CO2 levels to those predicted by your forecast from the ARIMA model in 1997 (i.e. “Task 3a”), we can see the results in our graph of True (Blue) vs Predicted (Green) CO2 Levels Over Time. Initially, there is a harmonious alignment between the true (blue) and predicted (green) CO2 levels, indicating that the ARIMA model effectively captured the underlying patterns in the earlier years. However, as time progresses, a noticeable trend emerges: the actual CO2 levels exhibit a more accelerated increase compared to what the ARIMA model predicted. This growing disparity suggests that there are evolving factors or trends influencing atmospheric CO2 concentrations that were not adequately accounted for in the original 1997 ARIMA model. These results are very similar to our previous graph vs our linear time model. Consistent with the broader Keeling Curve’s evolution from 1997 to the present, the yearly seasonal patterns and monthly variations in the true values stay consistent. This adherence to historical patterns underscores the persistent nature of the underlying dynamics of atmospheric CO2 concentrations. Furthermore, the observation that the predicted values fall behind the actual values over time hints at the potential influence of a quadratic term in the Keeling Curve. This suggests a more intricate relationship between time and CO2 levels than initially modeled.

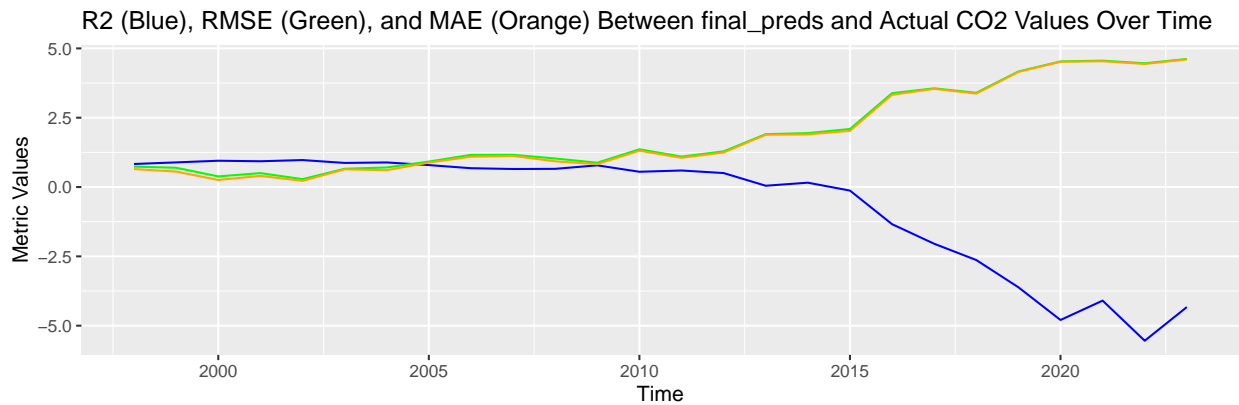
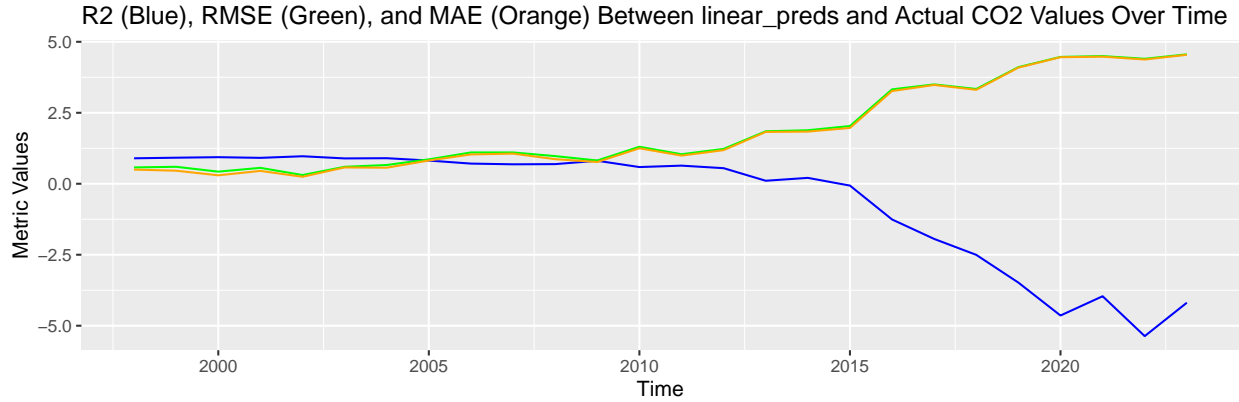


Evaluate the Performance of 1997 Linear and ARIMA Models

We initially predicted that atmospheric CO2 would cross 420ppm for the first time in Task 4a, but in this stage we calculated the truth and can see that it occurred in April 2022 vs our forecast of April 2024. Our models were close to the truth, being only two years off. This discrepancy suggests a lag in predicting the acceleration of CO2 levels, indicating the complexity of forecasting long-term environmental changes. Now we

continue to use the weekly data to generate a month-average series from 1997 to the present (month average series already generated during initial data ingestion), and compare the overall forecasting performance of our models from Parts 2a and 3b over the entire period.

In order to evaluate the performance of our Linear and ARIMA models, we take a look at their R2, RMSE, and MAE values over time. The first plot is the linear model performance and the second plot is the ARIMA model performance.



Train Best Models on Present Data

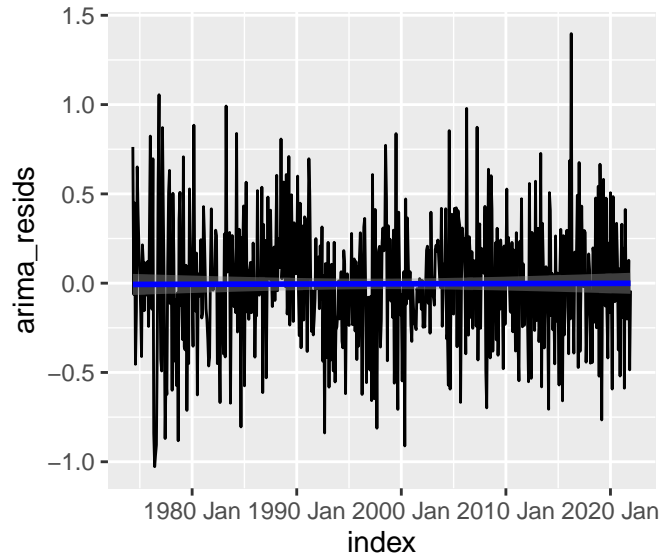
For training our best models on present data, we seasonally adjust the weekly NOAA data, and split both seasonally-adjusted (SA) and non-seasonally-adjusted (NSA) series into training and test sets, using the last two years of observations as the test sets, fitting ARIMA models for both SA and NSA series.

Our process outline was as follows: (if you want to see all of the output please see Appendix)

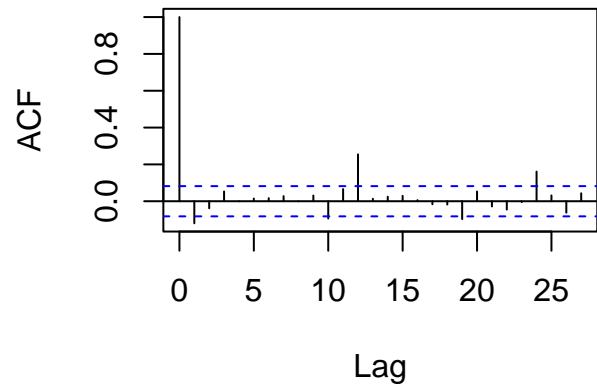
- We started by creating two versions of the CO2 data: one with seasonal adjustments (SA) and one without (NSA). These were then split into training and testing sets for model evaluation.
- The stationarity of the adjusted series (NSA and SA) was visually inspected, and the target variables over time were plotted.
- The autocorrelation and partial autocorrelation functions of the NSA and SA series were plotted to identify potential parameters for ARIMA modeling.
- ARIMA models were fitted to both the NSA and SA series using a grid search approach to identify optimal parameters. The models were evaluated in-sample and pseudo out-of-sample to measure their performance.
- A polynomial time-trend model was fitted to the seasonally-adjusted series, and its performance was compared to the ARIMA model.
- We then retrained the seasonally-adjusted series using a linear model for adjustment and then fit ARIMA models to the adjusted series.

- Fit the polynomial to seasonally adjusted data.
- For our final ARIMA model, we used linear model differencing because it allows us to reconstruct forecasts farther out into the future even when ARIMA model begins generating constant predictions.
- Evaluated residuals and performance on test set.
- Finally, we retried modeling non seasonally adjusted series using linear model to adjust and remade our final ARIMA model.

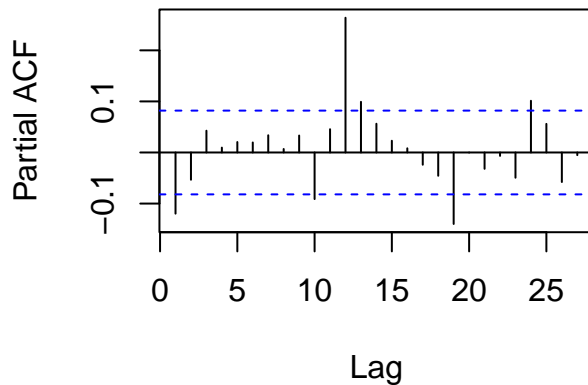
Below are visualizations used to evaluate our final ARIMA model.



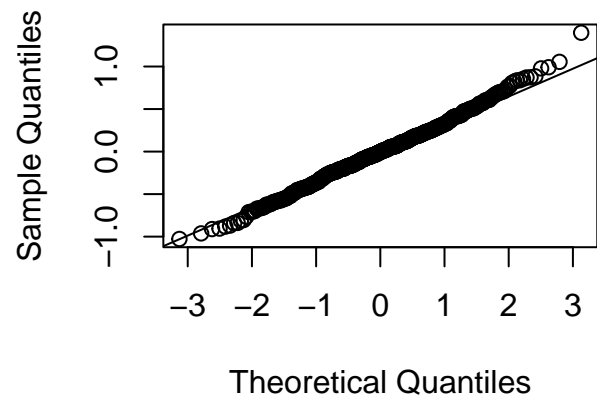
Series sa_train[["arima_resids"]]

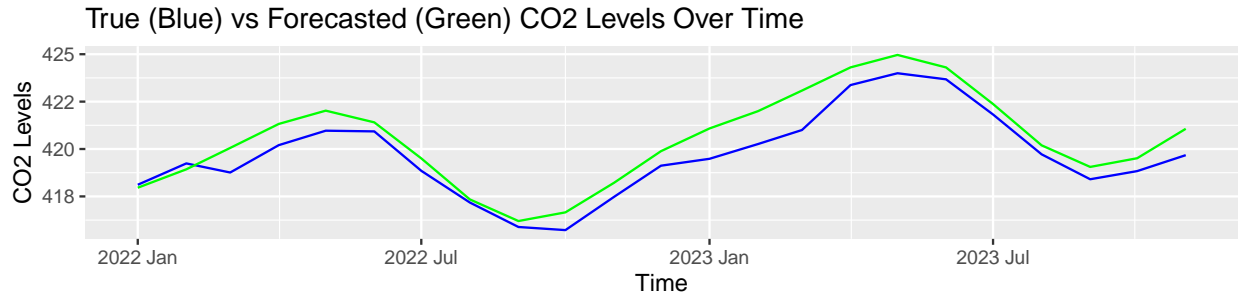


Series sa_train[["arima_resids"]]



Normal Q-Q Plot





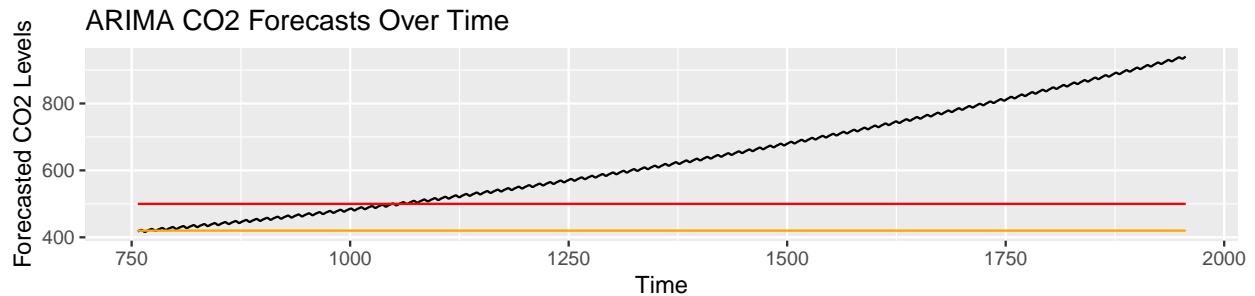
How Bad Could It Get? (Predictions of the Far Future)

Our code initiates by creating a non-seasonally adjusted subset (NSA) of CO2 data up to the year 2022. A linear model is then fitted to the NSA data, and its residuals are used to train an ARIMA model with parameters (5, 0, 4). We then generate future time points for CO2 prediction from the year 2022 to 2122. The features for the linear model are constructed, including indices, logarithmic transformations, and the month effect. We create a plot visualizing the ARIMA CO2 forecasts over time, including the forecasted CO2 levels, the threshold at 420, and the threshold at 500. And finally we identify the first and final times when the CO2 levels are predicted to cross the 420 ppm and 500ppm thresholds.

Using our non-seasonally adjusted data series, our generated predictions for when atmospheric CO2 is expected to be at 420 ppm and 500 ppm levels for the first and final times is as follows:

- First and Final Time at 420, Mar 2022 - Jul 2022
- First and Final time at 500, Apr 2046 - Jan 2047

Below is our prediction for atmospheric CO2 levels up to the year 2122, if the future keeps with the same pattern historically, we are fairly confident in these estimates.



Appendix

While our final results are reported here, in our complete notebook (Github Folder: Notebook) we examine alternative models and go into further assessment of the models. The purpose of this background information is to show more of the process in how we reached the conclusions shown above.