

project title*

project subtitle

a

Jiacan Sun

Kevin You

Ziheng Zhong

June 20, 2025

Table of contents

1	Introduction	2
2	Method	3
2.1	Data	3
2.2	Models	5
3	Result	9
3.1	Model 1	9
3.2	Model 2	12
4	Discussion	14
5	Appendix	14
5.1	Project Code	14
	References	15

*Code and data supporting this analysis is available at: [Link to repository](#).

1 Introduction

Climate change has become a critical global challenge, driven largely by rising carbon dioxide (CO_2) emissions from human activities. The Intergovernmental Panel on Climate Change (IPCC) warns that without immediate and deep emission reductions, the world is on track to exceed the Paris Agreement’s 1.5°C warming threshold, underscoring an urgent need to curb CO_2 from all sectors (Intergovernmental Panel on Climate Change (IPCC) 2022). In this context, the United States – historically the single largest national source of CO_2 – plays a pivotal role. Cumulatively, the U.S. has emitted roughly 25% of all fossil-fuel CO_2 since the industrial era began (Ritchie 2019), and it remains one of the top annual emitters today. Understanding the trajectory of U.S. emissions, especially from major sectors like transportation, industry, and energy generation, is therefore of global importance for climate change mitigation.

The motivation for this research is grounded in climate and energy policy relevance. The chosen sectors – transportation, industrial production, and power generation – are the dominant sources of U.S. greenhouse gas emissions (Keerthana et al. 2023). For example, in 2019 the transportation sector accounted for about 29% of U.S. CO_2 -equivalent emissions, followed by electricity generation (25%) and industrial processes (23%) (Keerthana et al. 2023). These activities collectively drive the bulk of national emissions, meaning any meaningful climate strategy must address each of them. Analyzing historical patterns in these sectors can reveal how past economic growth, technological changes, and policies (such as vehicle efficiency standards or power plant regulations) have impacted emissions. Moreover, forecasting future emissions is crucial for gauging progress toward sustainability goals. The U.S. has set ambitious targets under the Paris Agreement – pledging a 50–52% reduction in greenhouse gases by 2030 (from 2005 levels) and net-zero emissions by 2050 (Keerthana et al. 2023) – which heightens the real-world significance of this study. By projecting emissions trajectories, we can assess whether current trends align with these climate goals or if additional policy interventions may be required.

This paper will analyze monthly data from the U.S. Energy Information Administration (U.S. Energy Information Administration 2025), focusing on trends in energy consumption and production across major sectors. Our response variable is total CO_2 emissions, modeled as a function of six key predictors: transportation petroleum consumption, fossil fuel production, renewable energy generation, and energy consumption by the commercial, industrial, and residential sectors. Given the temporal nature of the data, we employ an ARIMA model and XXX, which will be discussed in more detail later. These two modeling approach allows us to assess both the cross-sectional and time-dependent dynamics influencing emissions, also enabling us to do forecasting .

The paper is structured as follows: Section 2 TBD, Section 3 TBD, Section 4 TBD.

2 Method

2.1 Data

The dataset is from the U.S. Energy Information Administration’s (EIA) Monthly Energy Review, which issues detailed energy statistics, including production, consumption, trade, and emissions metrics. The Monthly Energy Review is released on the last workday of each month. It contains preliminary data that are subject to revision in subsequent releases and historical data that may be updated when source publications are revised. For this project, data were extracted for the period ranging from January 1973 to February 2025, with a total of 626 records at a monthly frequency. The EIA collects these data through standardized survey forms completed by energy industry participants, following review protocols approved every three years to ensure consistency and quality across reporting entities. The dataset (U.S. Energy Information Administration 2025) is sourced from the the U.S. Energy Information Administration’s (EIA) Monthly Energy Review that issues detailed energy statistics, including production, consumption, trade, and emissions metrics. The Monthly Energy Review is released on the last workday of each month and contains preliminary data that are subject to revision in subsequent releases, as well as historical data that may be updated when source publications are revised. For this project, data were extracted for the period ranging from January 1973 to February 2025, with a total of 626 records at a monthly frequency. The EIA collects these data through standardized survey forms completed by energy industry participants, following review protocols approved every three years to ensure consistency and quality across reporting entities.

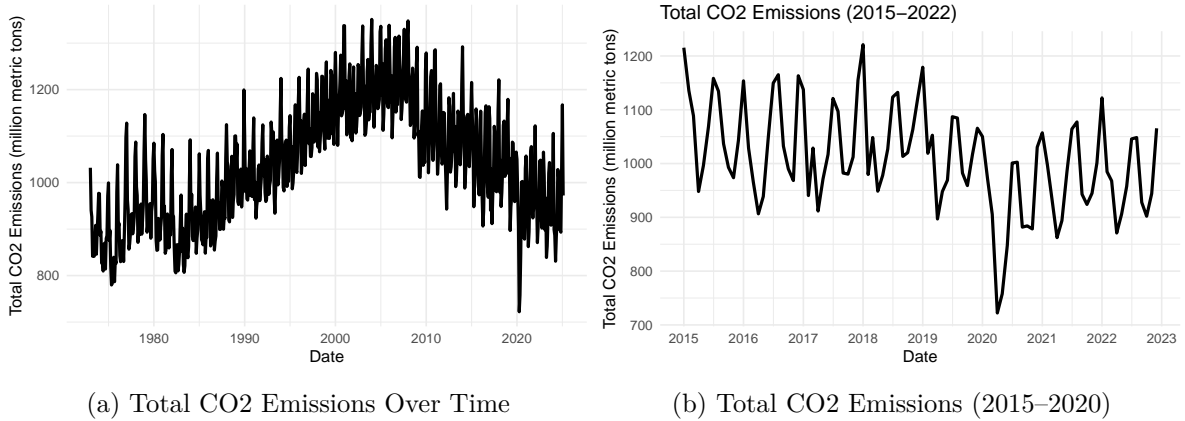


Figure 1: Data Overview

The cleaned dataset is constructed by concatenating monthly U.S. data from different categories. It includes energy production, sectoral energy use, transportation fuel use, and CO_2 emissions. Emissions from other sources are aggregated by date, and the total emissions across

all sectors are calculated. Additional statistics on petroleum consumption, fossil fuel and renewable energy output, and residential, commercial, and industrial energy consumption are merged after standardized Year and Month variables have been extracted. All datasets are joined in year and month.

The resulting dataset captures key indicators such as total CO_2 emissions (million metric tons), transportation petroleum consumption (million barrels per day), total fossil fuel and renewable energy output (quadrillion Btu), and sectoral energy consumption (trillion Btu). Units adhere to the conventions of the Monthly Energy Review, with energy metrics predominantly expressed in U.S. customary units—such as Btu for energy and barrels for petroleum volumes—and emissions in metric units.

From plot 1, we can say that total CO_2 emissions have a clear upward trend from the early 1970s through 2005 and then a decay from 2005 to the present. The sharp dip around 2020 likely reflects the economic slowdown during the COVID-19 pandemic. During this period, widespread lockdowns resulted in a significant reduction in transportation, industrial activity, and fossil fuel consumption, which contributed to the observed decline in emissions. The ACF plot shows that the total CO_2 emission is non-stationary since the lag-1 autocorrelation is significant and the ACF decays very slowly.

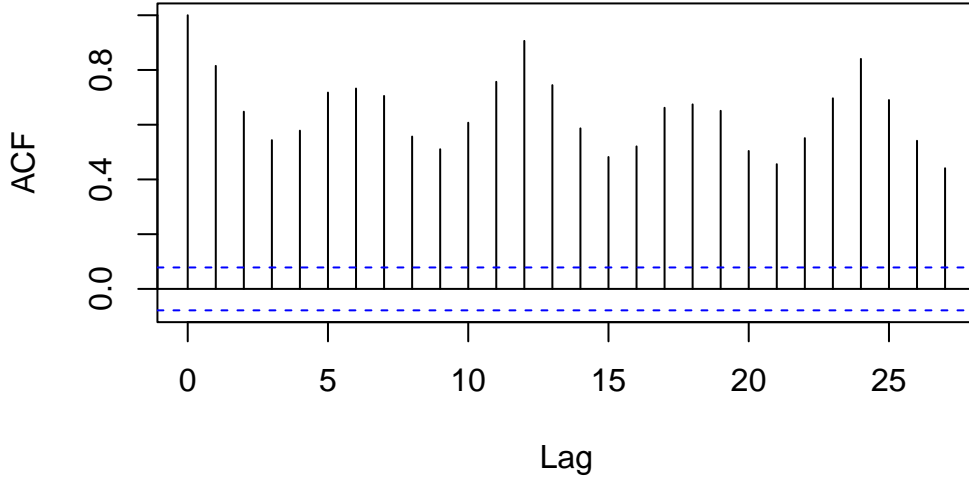


Figure 2: ACF of Total CO_2 Emissions

Key variables included in the cleaned dataset are total CO_2 emissions (million metric tons), transportation petroleum consumption (million barrels per day), total fossil fuels production

and total renewable energy production (quadrillion British thermal units [Btu]), and sectoral energy consumption in the commercial, industrial, and residential sectors (trillion Btu). Units adhere to the conventions of the Monthly Energy Review, with energy metrics predominantly expressed in U.S. customary units—such as Btu for energy and barrels for petroleum volumes—and emissions in metric units.

From plot 1, we can say that from the early 1970s through 2005, total CO_2 emissions has a clear upward trend, and then a decay from 2005 to present. The sharp dip around 2020 likely reflects the economic slowdown during the COVID-19 pandemic, everyone is in quarantine, hence there is less traffic and fossil fuel usage, which could be the reason of the dip.

By observing the ACF plot, we can conclude that the total CO_2 emission is non-stationary since the lag-1 autocorrelation is essentially 1.0, and the ACF decays very slowly, confirming the presence of non-stationary mean in the series.

2.2 Models

To model the relationship between total CO_2 emissions and the set of covariates we picked:

y_t - Total CO_2 emissions

z_{t1} - Transportation Petroleum Consumption

z_{t2} - Total Fossil Fuels Production

z_{t3} - Total Renewable Energy Production

z_{t4} - Commercial_Consumption

z_{t5} - Industrial_Consumption

z_{t6} - Residential_Consumption

while accounting for serial dependence in the errors, we followed the four-step procedure for regression with auto-correlated errors:

First, we fit an ordinary regression of y_t on z_{t1}, \dots, z_{t6} (acting as if the errors are uncorrelated).

$$y_t = \sum_{j=1}^6 \beta_j z_{tj} + x_t$$

From this fit we can retain the residuals by

$$\hat{x}_t = y_t - \sum_{j=1}^r \beta_j z_{tj}$$

After fitting the initial linear regression model, we examined the residuals to check for autocorrelation using autocorrelation (ACF) and partial autocorrelation (PACF) plots, as well as the Box-Ljung test. These diagnostics showed strong and persistent correlation patterns, indicating that the residuals were not independent over time.

Table 1: Model 1 Summary

Variable	Estimate	Std. Error	t value	Pr(> t)	Significance
(Intercept)	294.0642	21.0351	13.9797	0.0000	***
Transportation	1.2353	0.2103	5.8735	0.0000	***
Fossil Fuel Prod.	-38.0689	4.3954	-8.6610	0.0000	***
Renewable Energy Prod.	-256.3311	26.5226	-9.6646	0.0000	***
Commercial	0.5023	0.0133	37.8640	0.0000	***
Industrial	0.1639	0.0081	20.1718	0.0000	***
Residential	0.0015	0.0058	0.2667	0.7898	

In our case, our model is

$$y_t = 294.1 + 1.235z_{t1} - 38.07z_{t2} - 256.3z_{t3} + 0.5023z_{t4} + 0.1639z_{t5} + 0.001546z_{t6} + x_t$$

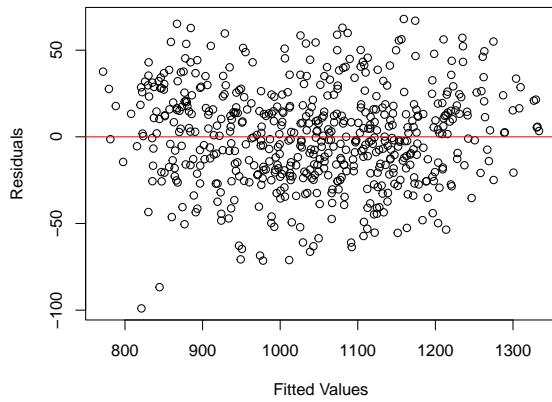
In the scatter plot of residuals against the fitted values (Figure 3a), there is no clear funnel or curvature, which suggests that the variance of the errors is roughly constant across the range of predicted CO_2 emissions.

However, when we plot the residuals in time order, a pronounced seasonal “saw tooth” pattern emerges—residuals swing positive and negative in a regular annual pattern (Figure 3b). This indicates that our purely cross-sectional predictors haven’t captured the strong yearly cycle in emissions, and so that seasonality remains in the error term.

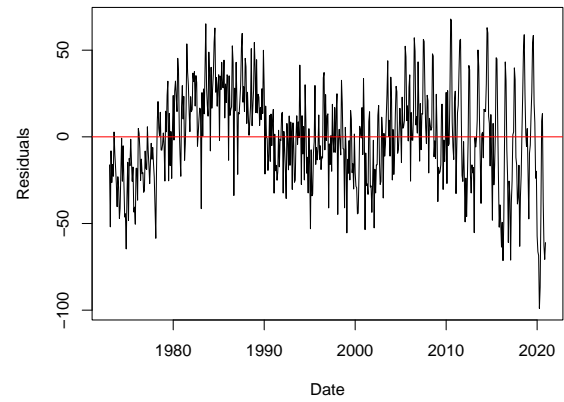
The histogram of residuals (Figure 3c) is roughly symmetric around zero and appears approximately bell-shaped, albeit with mildly heavier tails than a perfect normal distribution.

The Q-Q plot (Figure 3d) confirms this: most points lie very close to the 45° line, indicating near-normality, but the extreme ends (both lower and upper tails) deviate slightly. In practice, these small departures from Gaussianity are unlikely to invalidate inference, especially once we explicitly model the autocorrelation and seasonality.

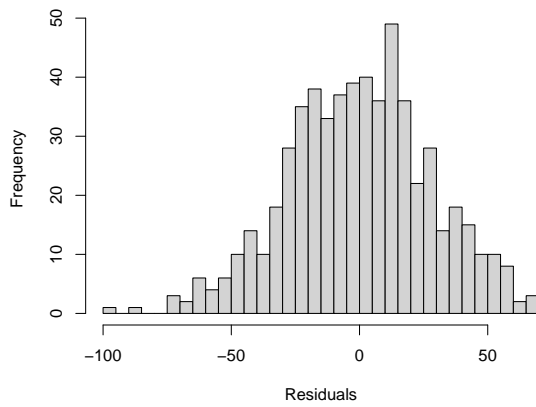
Then we examined the sample ACF and PACF of \hat{x}_t to identify candidate ARMA(p,q) structures. Potential orders p and q were chosen by the usual cut-off and tail-off patterns in these plots. Finally we applied `forecast::auto.arima()` to select an appropriate ARIMA(p,d,q) model



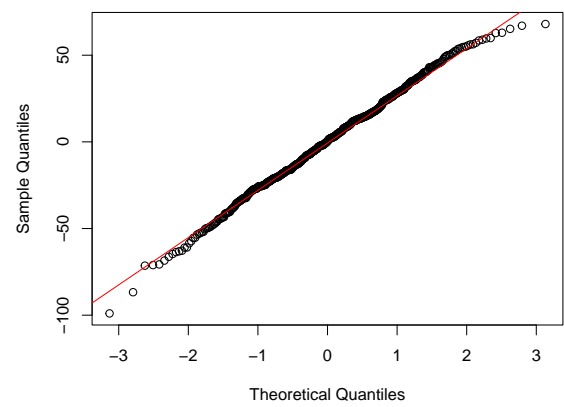
(a) Residuals vs Fitted



(b) Residuals Over Time



(c) Histogram of Residuals



(d) Normal Q-Q Plot

Figure 3: Model 1 Diagnostic Plots

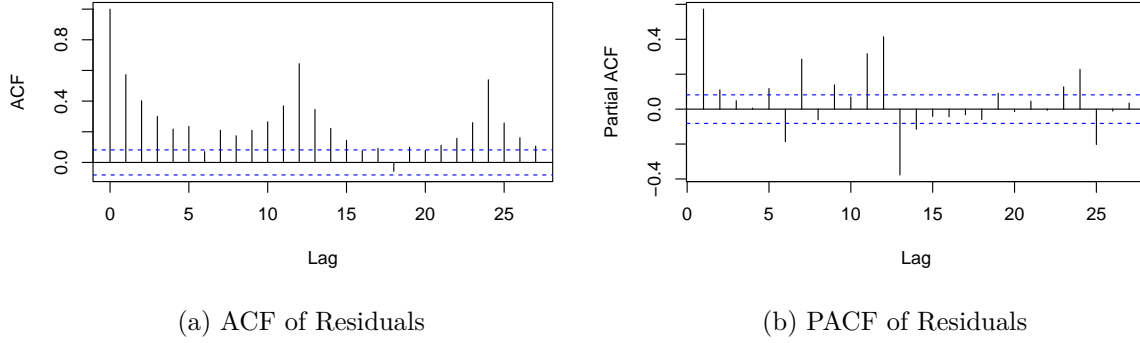


Figure 4: Model 1 ACF & PACF

Figure 4a shows the autocorrelation of residuals from our linear model. We observe significant spikes at multiple lags, especially early ones, which exceed the blue dashed significance bands. This means the residuals are not independent over time, and there is clear serial correlation remaining in the model. This suggests that our linear model has not fully captured the temporal structure of the data.

Figure 4b shows partial autocorrelations of the residuals. We see large spikes at the first few lags, indicating that past values have a direct influence on the current residuals, even after accounting for previous lags. This supports the idea that an autoregressive structure may still be present and needs to be modeled, possibly with an ARIMA or time series component.

Table 2: Box-Ljung Test on Model Residuals

	Test	Statistic	df	p.value
X-squared	Box-Ljung	967.69	20	<0.001

The Box-Ljung test (Table 2) checks whether residuals are independently distributed. Here, the test statistic is 967.69 with 20 degrees of freedom, and the p-value is less than 0.001. This result strongly rejects the null hypothesis of no autocorrelation, confirming that the residuals still contain time-dependent patterns and the current model is insufficient to capture them.

Table 3: ARIMA Model Summary on OLS Residuals

Metric	Value
Sigma ²	468.2533
Log Likelihood	-2581.0715
AIC	5178.1431

Table 3: ARIMA Model Summary on OLS Residuals

Metric	Value
AICc	5178.3975
BIC	5212.9780
ME	-0.0920
RMSE	21.4884
MAE	16.7951
MPE	93.5415
MAPE	256.5873
MASE	0.8007
ACF1	0.0280

To address issues mentioned above, we applied an ARIMA model to the residuals from our linear regression. We used the `auto.arima()` function to select the best-fitting ARIMA model for the residuals. More details and results are presented in [?@sec-res](#).

We also decided to explore a machine learning approach using a Long Short-Term Memory (LSTM) neural network. LSTM is well-suited for sequential data and can capture long-term temporal dependencies without requiring stationarity. The model was trained... To Be Added

Full implementation details and evaluation results are provided in [?@sec-res](#).

3 Result

3.1 Model 1

Continuing with Method section, our `auto.arima()` call picked an ARIMA(5,1,2) model on the OLS residuals \hat{x}_t . Concretely, letting $\nabla \hat{x}_t = \hat{x}_t - \hat{x}_{t-1}$, the fitted model is

$$(1 - 0.2594B + 0.3250B^2 + 0.0259B^3 - 0.1480B^4 + 0.1365B^5)(1 - B)\hat{x}_t = (1 - 0.2066B - 0.6959B^2)w_t$$

where B is the backshift operator and w_t is white noise with estimated $\sigma^2(w_t) = 468.3$. AIC of 5178.1 confirms this is the most parsimonious high-order fit `auto.arima` found, and the small training-set ACF1 (≈ 0.028) suggests remaining serial dependence is minimal after differencing and fitting these parameters.

Table 4: ARIMACoefficients and Standard Errors

	Term	Estimate	Std_Error
ar1	ar1	-0.2594	0.0618

Table 4: ARIMACoefficients and Standard Errors

	Term	Estimate	Std_Error
ar2	ar2	0.3250	0.0494
ar3	ar3	0.0259	0.0466
ar4	ar4	-0.1480	0.0456
ar5	ar5	0.1365	0.0465
ma1	ma1	-0.2066	0.0473
ma2	ma2	-0.6959	0.0418

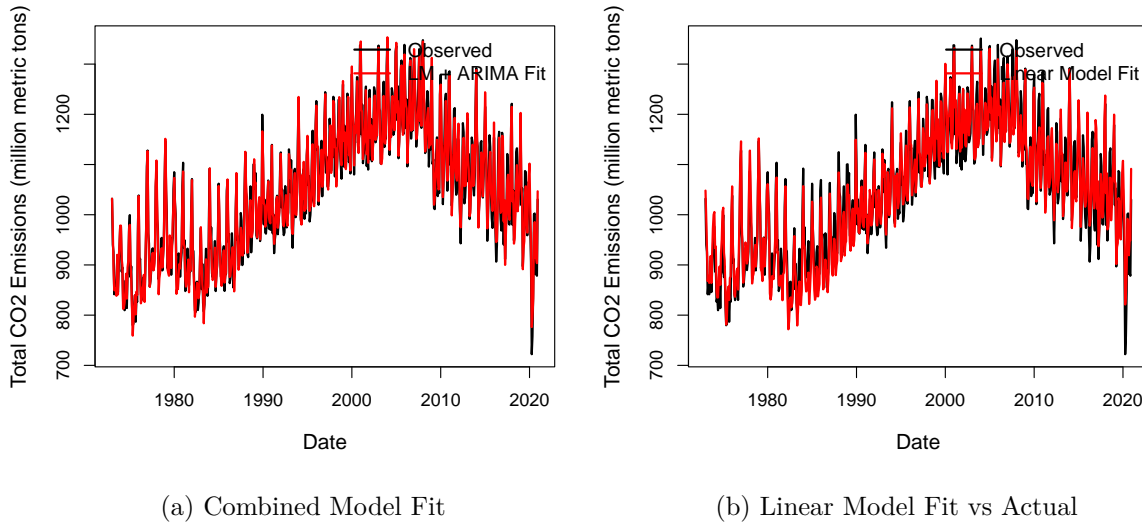


Figure 5: Model Fittings

Figure 5 compares the performance of two models in capturing U.S. CO_2 emissions over time. In both plots, the black line represents the observed emissions, while the red line shows the model predictions. Figure 5a shows the combined model (linear + ARIMA), which closely follows the seasonal and long-term trends, especially in the mid and later periods. Figure 5b displays the linear model alone, which captures the overall trend but misses the recurring seasonal spikes and dips.

Table 5: Comparison of Linear Model vs. Linear + ARIMA Model

Model	RSS	RMSE
Linear Model	451898.4	28.0097
Linear + ARIMA Model	265967.9	21.4884

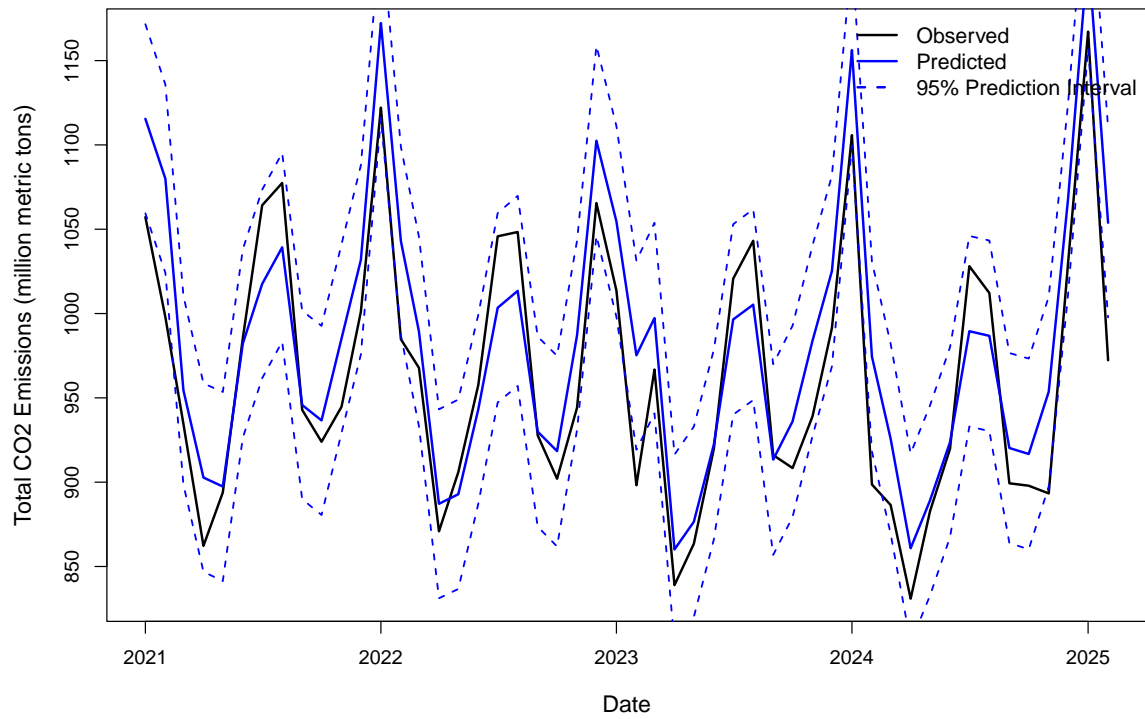


Figure 6: RSS and RMSE Comparison of Linear vs. Linear + ARIMA Models

Figure 6 shows how well the combined Linear + ARIMA model predicts CO_2 emissions for the years 2021 to 2025. The black line represents the actual observed emissions, while the solid blue line shows the model's predicted values, and the dashed blue lines indicate the 95% prediction interval. The predicted values follow the seasonal ups and downs of the observed data quite closely, staying mostly within the prediction bounds. This suggests that the model is effective at capturing both the trend and seasonal patterns in the data.

3.2 Model 2

Figure 7: Loss Curve

Figure 7 shows how the model's error changes over time during training. Specifically, it plots the Mean Squared Error (MSE) for both the training and validation data across 40 epochs. At the beginning, both training and validation losses are relatively high, especially the training loss. As training progresses, both losses gradually decrease. This means the model is learning and improving its accuracy over time. Around the middle of the training, the validation loss starts to follow the same trend as the training loss and eventually becomes even lower. This is a good sign, showing that the model is not overfitting and is able to generalize well to unseen data. Overall, the loss curve shows that the model training was successful and well-regularized.

Figure 8 compares the model's predicted CO_2 emissions (in orange) with the actual true emissions (in blue) for the years after 2022. The two lines follow very similar patterns, showing

Figure 8: Test Forecast

that the model has learned the seasonal trends in the data. Although the predicted values are not exactly the same as the true values, they are close and show the same ups and downs. This suggests the model is able to capture important patterns over time, like seasonal cycles. The small differences between the predicted and actual values are expected and acceptable. In general, the forecast shows that the LSTM model performs well in predicting future CO_2 emissions.

4 Discussion

5 Appendix

5.1 Project Code

References

- Intergovernmental Panel on Climate Change (IPCC). 2022. “Climate Change 2022: Mitigation of Climate Change.”
- Keerthana, Krishnamurthy Baskar, Shih-Wei Wu, Mu-En Wu, and Thangavelu Kokulnathan. 2023. “The United States Energy Consumption and Carbon Dioxide Emissions: A Comprehensive Forecast Using a Regression Model.” *Sustainability* 15 (10): 7932. <https://doi.org/10.3390/su15107932>.
- Ritchie, Hannah. 2019. “Who Has Contributed Most to Global CO2 Emissions?” *Our World in Data*.
- U.S. Energy Information Administration. 2025. “Monthly Energy Review.” U.S. Energy Information Administration; <https://www.eia.gov/totalenergy/data/monthly/>.