

Featured Models: RF, XGBoost, SARIMAX, OLS

Key Concepts: Diagnostics, Seasonality, Humidity Regimes

Modeling Pillars & Behavioral Themes

Toolkit Behind the Analysis: *Python, R, Power BI, Gamma App, Kaggle GitHub*

Source: "Steel Industry Energy and Emissions Data," available on Kaggle

## Decoding Steel Plant Systems: Energy Use, Emissions, and Optimization

Optimizing Energy Use Reducing CO<sub>2</sub> Output at a Steel plant: A Data-Driven Approach

# Summary & Portfolio Overview

This portfolio delivers a layered understanding of energy demand and emissions behavior at a steel plant—blending environmental diagnostics, operational modeling, and temporal forecasting.

Through exploratory analysis, reactive power metrics, humidity regimes, and seasonality, I uncovered nonlinear relationships between usage and emissions.

Modeling spanned OLS, Random Forest, XGBoost, SARIMAX, and ARIMA; each selected and refined with careful diagnostic reflection.

Beyond modeling, this journey included hands-on development in Python, R, Power BI, Excel, Kaggle, and GitHub; building not just predictions, but pipelines, visualizations, and narratives that reveal how energy systems breathe and behave.

*Explore deeper sections to see: feature engineering breakthroughs, seasonality decomposition, imputation trade-offs, and diagnostic plots that uncover system "personality."*

# Portfolio Overview

## Sections at a Glance



### Data Collection & Transformation

Kaggle dataset + weather merge, datetime alignment, CO<sub>2</sub> imputation, scaling



### Feature Engineering & Diagnostics

Humidity bins, load rank, reactive power signals, outlier treatment



### Temporal Forecasting

SARIMAX, ARIMA — autocorrelation, seasonal cycles, temperature sensitivity



### Cost Impact & Operational Value

Quantified savings (\$7,505+), emission insights, optimization pathways



### Exploratory Analysis & Insights

Distributions, correlations, regime detection, weekday effects, seasonality



### Structural Modeling

OLS, RF, XGBoost — residual analysis, Q-Q plots, regime comparisons



### Behavioral & Rhythm Models

log\_kWh model, weekday/weekend shifts, Tuesday spikes



### Final Reflections & Conclusion

Modeling arc, personal growth, data storytelling, next steps

## **Section 1: Data Collection and Transformation**

# Data Preprocessing Challenges & Solutions

## DATA WRANGLING

1. Grabbing the dataset from Kaggle ([Steel\\_industry\\_data.csv\(2.73 MB\)](#))
2. Fixed character encoding issues (converted "Â°C" to "°C")
  - a. There are times I still like to revisit the raw data and transform differently
3. Researching and finding Weather Data set from the same time
  - a. The dataset was about 30 miles away from the Steelplant
4. Resolved datetime format mismatches between energy and weather datasets
  - a. Resampled weather data to align with energy consumption timestamps
5. Merged datasets to create a comprehensive analytical base
6. Implemented strategic handling of missing values
  - a. Continuous transforming for strategic modeling

```
humidity_col = 'relative_humidity_2m (%)'  
temperature_col = 'temperature_2m (°C)'  
CO2_col = 'CO2(tCO2)'
```

date	day	day_of_week	day_name	week	month	month_name	year
0 2018-01-01 00:15:00	1	0	Monday	1	1	January	2018
1 2018-01-01 00:30:00	1	0	Monday	1	1	January	2018
2 2018-01-01 00:45:00	1	0	Monday	1	1	January	2018
3 2018-01-01 01:00:00	1	0	Monday	1	1	January	2018
4 2018-01-01 01:15:00	1	0	Monday	1	1	January	2018

*Note: Although painful, after all my analysis and modeling, I am very pleased that I took the time to find the corresponding weather data.*

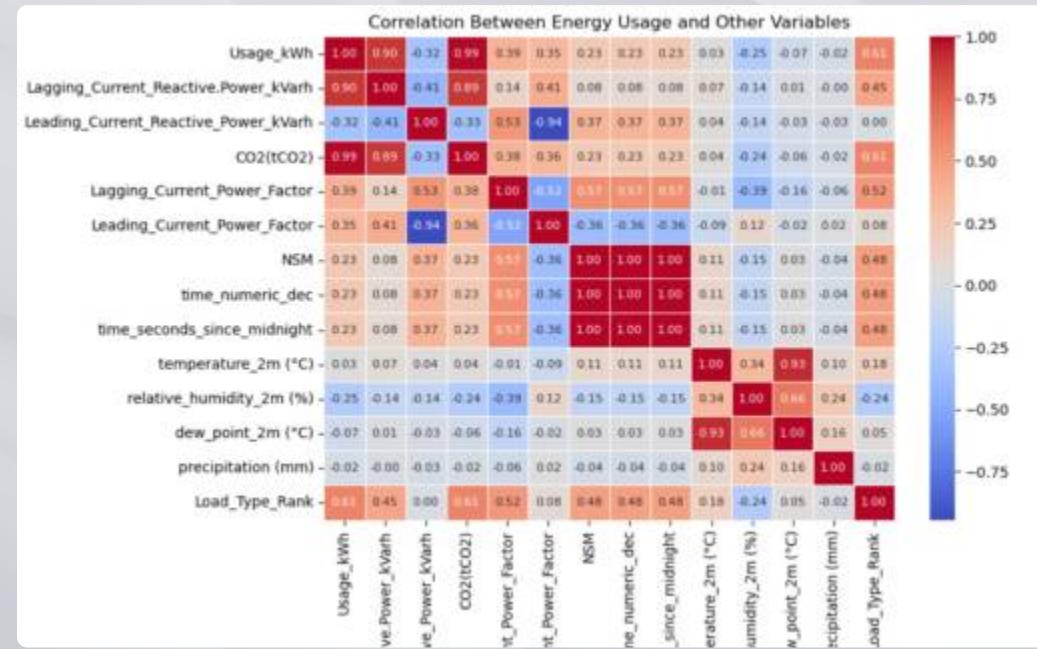
date	day	day_of_week	day_name	week	month	month_name	year	consumption_twh	temp_min	temp_max	humidity_min	humidity_max	precipitation_mm	cloud_type	weather
0 2018-01-01 00:15:00	1	0	Monday	1	1	January	2018	34400	68.2	71	82	9	0	P	Sunny
1 2018-01-01 00:30:00	1	0	Monday	1	1	January	2018	34400	68.2	70	88	52	0	P	Sunny
2 2018-01-01 00:45:00	1	0	Monday	1	1	January	2018	24400	68.2	69	88	52	0	P	Sunny
3 2018-01-01 01:00:00	1	0	Monday	1	1	January	2018	25000	68.2	68	88	52	0	P	Sunny
4 2018-01-01 01:15:00	1	0	Monday	1	1	January	2018	25000	68.2	67	88	52	0	P	Sunny
5 2018-01-01 01:30:00	1	0	Monday	1	1	January	2018	25000	68.2	66	88	52	0	P	Sunny
6 2018-01-01 01:45:00	1	0	Monday	1	1	January	2018	25000	68.2	65	88	52	0	P	Sunny
7 2018-01-01 02:00:00	1	0	Monday	1	1	January	2018	25000	68.2	64	88	52	0	P	Sunny
8 2018-01-01 02:15:00	1	0	Monday	1	1	January	2018	25000	68.2	63	88	52	0	P	Sunny
9 2018-01-01 02:30:00	1	0	Monday	1	1	January	2018	25000	68.2	62	88	52	0	P	Sunny
10 2018-01-01 02:45:00	1	0	Monday	1	1	January	2018	25000	68.2	61	88	52	0	P	Sunny
11 2018-01-01 03:00:00	1	0	Monday	1	1	January	2018	25000	68.2	60	88	52	0	P	Sunny
12 2018-01-01 03:15:00	1	0	Monday	1	1	January	2018	25000	68.2	59	88	52	0	P	Sunny
13 2018-01-01 03:30:00	1	0	Monday	1	1	January	2018	25000	68.2	58	88	52	0	P	Sunny
14 2018-01-01 03:45:00	1	0	Monday	1	1	January	2018	25000	68.2	57	88	52	0	P	Sunny
15 2018-01-01 04:00:00	1	0	Monday	1	1	January	2018	25000	68.2	56	88	52	0	P	Sunny
16 2018-01-01 04:15:00	1	0	Monday	1	1	January	2018	25000	68.2	55	88	52	0	P	Sunny
17 2018-01-01 04:30:00	1	0	Monday	1	1	January	2018	25000	68.2	54	88	52	0	P	Sunny
18 2018-01-01 04:45:00	1	0	Monday	1	1	January	2018	25000	68.2	53	88	52	0	P	Sunny
19 2018-01-01 05:00:00	1	0	Monday	1	1	January	2018	25000	68.2	52	88	52	0	P	Sunny
20 2018-01-01 05:15:00	1	0	Monday	1	1	January	2018	25000	68.2	51	88	52	0	P	Sunny
21 2018-01-01 05:30:00	1	0	Monday	1	1	January	2018	25000	68.2	50	88	52	0	P	Sunny
22 2018-01-01 05:45:00	1	0	Monday	1	1	January	2018	25000	68.2	49	88	52	0	P	Sunny
23 2018-01-01 06:00:00	1	0	Monday	1	1	January	2018	25000	68.2	48	88	52	0	P	Sunny
24 2018-01-01 06:15:00	1	0	Monday	1	1	January	2018	25000	68.2	47	88	52	0	P	Sunny
25 2018-01-01 06:30:00	1	0	Monday	1	1	January	2018	25000	68.2	46	88	52	0	P	Sunny
26 2018-01-01 06:45:00	1	0	Monday	1	1	January	2018	25000	68.2	45	88	52	0	P	Sunny
27 2018-01-01 07:00:00	1	0	Monday	1	1	January	2018	25000	68.2	44	88	52	0	P	Sunny
28 2018-01-01 07:15:00	1	0	Monday	1	1	January	2018	25000	68.2	43	88	52	0	P	Sunny
29 2018-01-01 07:30:00	1	0	Monday	1	1	January	2018	25000	68.2	42	88	52	0	P	Sunny
30 2018-01-01 07:45:00	1	0	Monday	1	1	January	2018	25000	68.2	41	88	52	0	P	Sunny
31 2018-01-01 08:00:00	1	0	Monday	1	1	January	2018	25000	68.2	40	88	52	0	P	Sunny
32 2018-01-01 08:15:00	1	0	Monday	1	1	January	2018	25000	68.2	39	88	52	0	P	Sunny
33 2018-01-01 08:30:00	1	0	Monday	1	1	January	2018	25000	68.2	38	88	52	0	P	Sunny
34 2018-01-01 08:45:00	1	0	Monday	1	1	January	2018	25000	68.2	37	88	52	0	P	Sunny
35 2018-01-01 09:00:00	1	0	Monday	1	1	January	2018	25000	68.2	36	88	52	0	P	Sunny
36 2018-01-01 09:15:00	1	0	Monday	1	1	January	2018	25000	68.2	35	88	52	0	P	Sunny
37 2018-01-01 09:30:00	1	0	Monday	1	1	January	2018	25000	68.2	34	88	52	0	P	Sunny
38 2018-01-01 09:45:00	1	0	Monday	1	1	January	2018	25000	68.2	33	88	52	0	P	Sunny
39 2018-01-01 09:59:00	1	0	Monday	1	1	January	2018	25000	68.2	32	88	52	0	P	Sunny
40 2018-01-01 10:14:00	1	0	Monday	1	1	January	2018	25000	68.2	31	88	52	0	P	Sunny
41 2018-01-01 10:29:00	1	0	Monday	1	1	January	2018	25000	68.2	30	88	52	0	P	Sunny
42 2018-01-01 10:44:00	1	0	Monday	1	1	January	2018	25000	68.2	29	88	52	0	P	Sunny
43 2018-01-01 10:59:00	1	0	Monday	1	1	January	2018	25000	68.2	28	88	52	0	P	Sunny
44 2018-01-01 11:14:00	1	0	Monday	1	1	January	2018	25000	68.2	27	88	52	0	P	Sunny
45 2018-01-01 11:29:00	1	0	Monday	1	1	January	2018	25000	68.2	26	88	52	0	P	Sunny
46 2018-01-01 11:44:00	1	0	Monday	1	1	January	2018	25000	68.2	25	88	52	0	P	Sunny
47 2018-01-01 11:59:00	1	0	Monday	1	1	January	2018	25000	68.2	24	88	52	0	P	Sunny
48 2018-01-01 12:14:00	1	0	Monday	1	1	January	2018	25000	68.2	23	88	52	0	P	Sunny
49 2018-01-01 12:29:00	1	0	Monday	1	1	January	2018	25000	68.2	22	88	52	0	P	Sunny
50 2018-01-01 12:44:00	1	0	Monday	1	1	January	2018	25000	68.2	21	88	52	0	P	Sunny
51 2018-01-01 12:59:00	1	0	Monday	1	1	January	2018	25000	68.2	20	88	52	0	P	Sunny
52 2018-01-01 13:14:00	1	0	Monday	1	1	January	2018	25000	68.2	19	88	52	0	P	Sunny
53 2018-01-01 13:29:00	1	0	Monday	1	1	January	2018	25000	68.2	18	88	52	0	P	Sunny
54 2018-01-01 13:44:00	1	0	Monday	1	1	January	2018	25000	68.2	17	88	52	0	P	Sunny
55 2018-01-01 13:59:00	1	0	Monday	1	1	January	2018	25000	68.2	16	88	52	0	P	Sunny
56 2018-01-01 14:14:00	1	0	Monday	1	1	January	2018	25000	68.2	15	88	52	0	P	Sunny
57 2018-01-01 14:29:00	1	0	Monday	1	1	January	2018	25000	68.2	14	88	52	0	P	Sunny
58 2018-01-01 14:44:00	1	0	Monday	1	1	January	2018	25000	68.2	13	88	52	0	P	Sunny
59 2018-01-01 14:59:00	1	0	Monday	1	1	January	2018	25000	68.2	12	88	52	0	P	Sunny
60 2018-01-01 15:14:00	1	0	Monday	1	1	January	2018	25000	68.2	11	88	52	0	P	Sunny
61 2018-01-01 15:29:00	1	0	Monday	1	1	January	2018	25000	68.2	10	88	52	0	P	Sunny
62 2018-01-01 15:44:00	1	0	Monday	1	1	January	2018	25000	68.2	9	88	52	0	P	Sunny
63 2018-01-01 15:59:00	1	0	Monday	1	1	January	2018	25000	68.2	8	88	52	0	P	Sunny
64 2018-01-01 16:14:00	1	0	Monday	1	1	January	2018	25000	68.2	7	88	52	0	P	Sunny
65 2018-01-01 16:29:00	1	0	Monday	1	1	January	2018	25000	68.2	6	88	52	0	P	Sunny
66 2018-01-01 16:44:00	1	0	Monday	1	1	January	2018	25000	68.2	5	88	52	0	P	Sunny
67 2018-01-01 16:59:00	1	0	Monday	1	1	January	2018	25000	68.2	4	88	52	0	P	Sunny
68 2018-01-01 17:14:00	1	0	Monday	1	1	January	2018	25000	68.2	3	88	52	0	P	Sunny
69 2018-01-01 17:29:00	1	0	Monday	1	1	January	2018	25000	68.2	2	88	52	0	P	Sunny
70 2018-01-01 17:44:00	1	0	Monday	1	1	January	2018	25000	68.2	1	88	52	0	P	Sunny
71 2018-01-01 17:59:00	1	0	Monday	1	1	January	2018	25000	68.2	0	88	52	0	P	Sunny
72 2018-01-01 18:14:00	1	0	Monday	1	1	January	2018	25000	68.2	-1	88	52	0	P	Sunny
73 2018-01-01 18:29:00	1	0	Monday	1	1	January	2018	25000	68.2	-2	88	52	0	P	Sunny
74 2018-01-01 18:44:00	1	0	Monday	1	1	January	2018								

## **Section 2: Exploratory Analysis and Insights**

# Exploratory and Data Insights

## Getting to know the Data

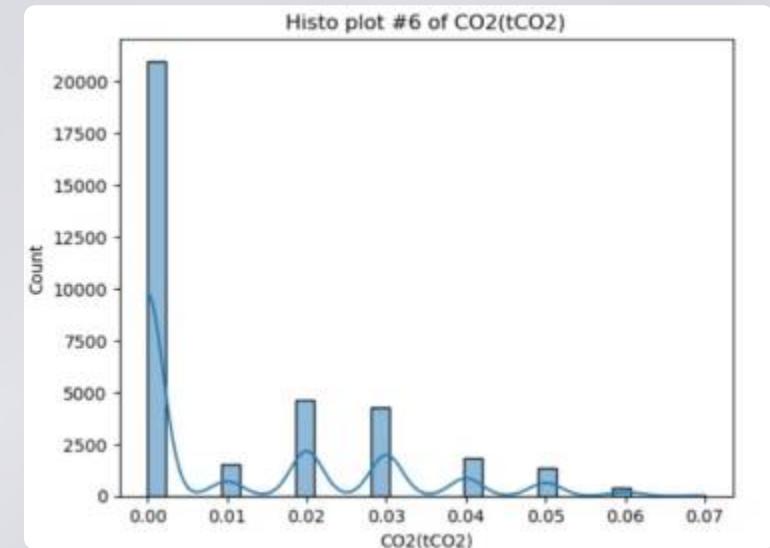
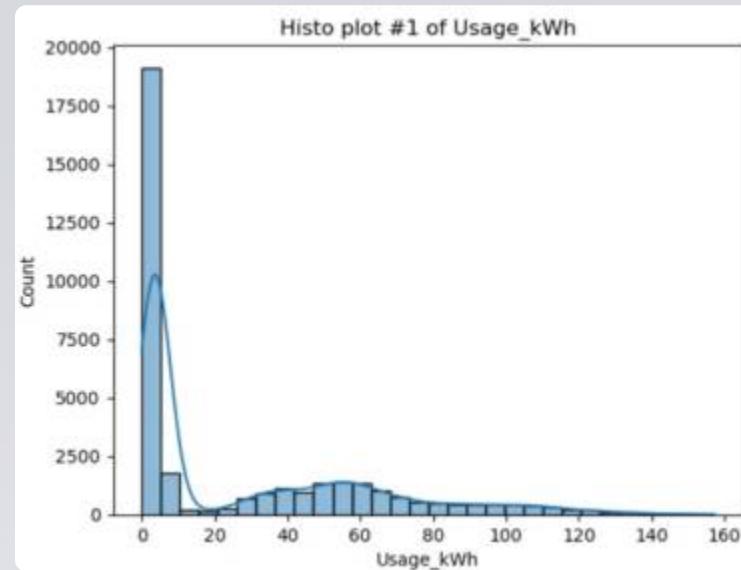
Before diving into modeling, I always take time to explore and interact with the data, not just to clean it, but to understand it. This includes visualizing distributions, examining trends over time, and checking for correlations or unexpected patterns. These early insights often spark better feature design and sharper hypotheses.



- Initial heatmap analysis revealed a strong positive correlation between energy usage (kWh) and CO<sub>2</sub> emissions; a relationship consistent with expected physical output, where increased energy production corresponds with higher emissions. This was key to guiding modeling decisions.

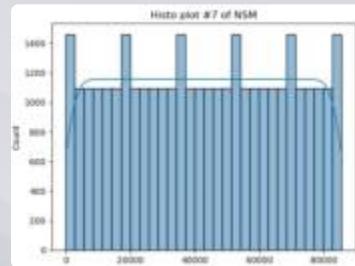
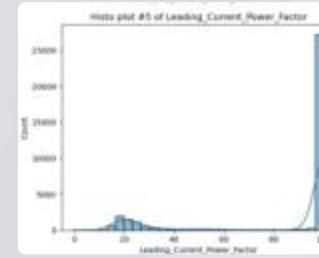
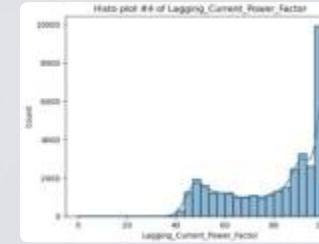
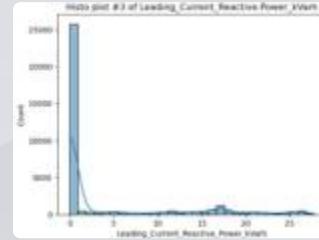
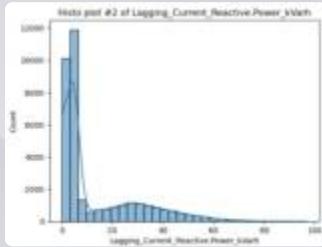
# Exploratory Data Insights - Initial Distributions

As someone trained in statistics, I believe it's essential to understand the shape and structure of the data before modeling. This includes assessing normality, identifying skewness or multimodality, and applying transformations where needed. These early checks help ensure that later models are both appropriate and interpretable.

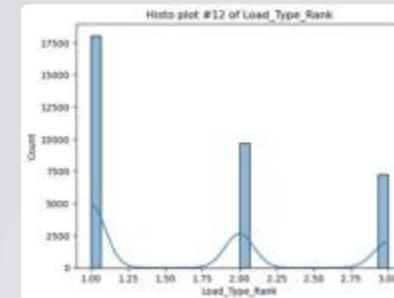


In this case, histograms of CO<sub>2</sub> and kWh revealed skewed distributions, leading to the application of transformation (log) on energy usage and an imputation strategy for zero-inflated CO<sub>2</sub> values.

# Exploratory Data Insights: Energy Distributions

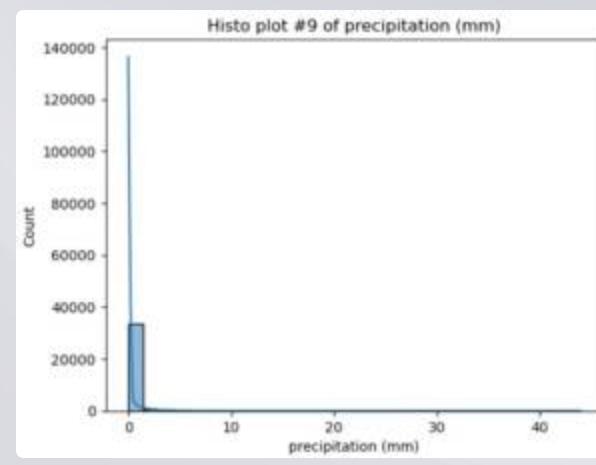
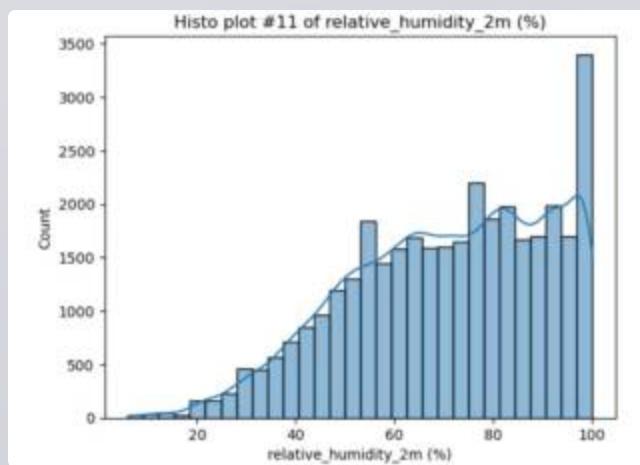
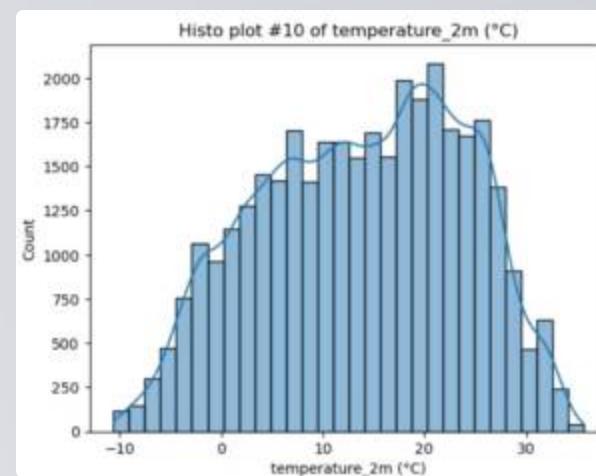
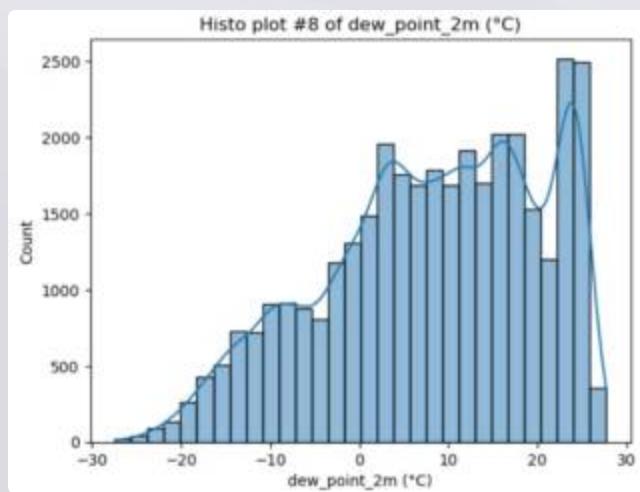


While several energy-related variables exhibited skewed distributions with some leaning heavily toward higher values. I also explored features like NSM and Rank, which are inherently nonparametric in nature. > > As someone with a background in statistics, I find it important not only to identify when a transformation (e.g., log scaling) might improve model performance, but also to understand the underlying shape of the distribution itself. These early diagnostics help determine whether a parametric or nonparametric approach is more appropriate.



Time series decomposition revealed seasonal components that align with production schedules. After log-transformation, the relationship between variables became more linear, supporting our modeling approach. Weather variables showed moderate influence on energy efficiency.

# Exploratory Data Insights: Environmental Distributions

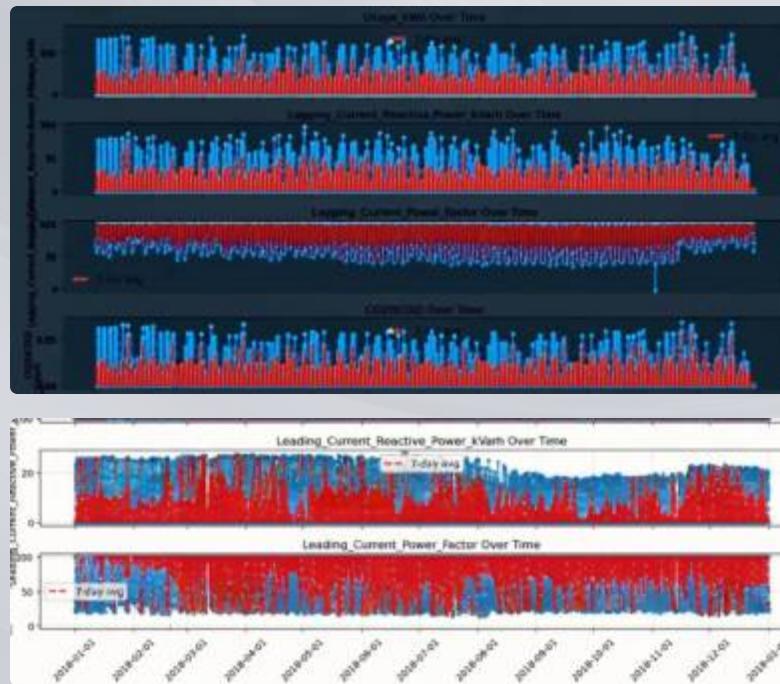


Histogram analysis showed that temperature and dew point followed approximately normal distributions, supporting their direct use in modeling without transformation. In contrast, humidity exhibited a strong right skew, suggesting the site regularly experiences high-humidity conditions. This insight guided the decision to engineer humidity-based bins and explore its role in emission variability. > > Another feature appeared to have limited variation and predictive value and was subsequently excluded from modeling to preserve parsimony. Precipitation appears to have limited data except because of the skew to zero won't be a main variable for the modeling but will still include initially.

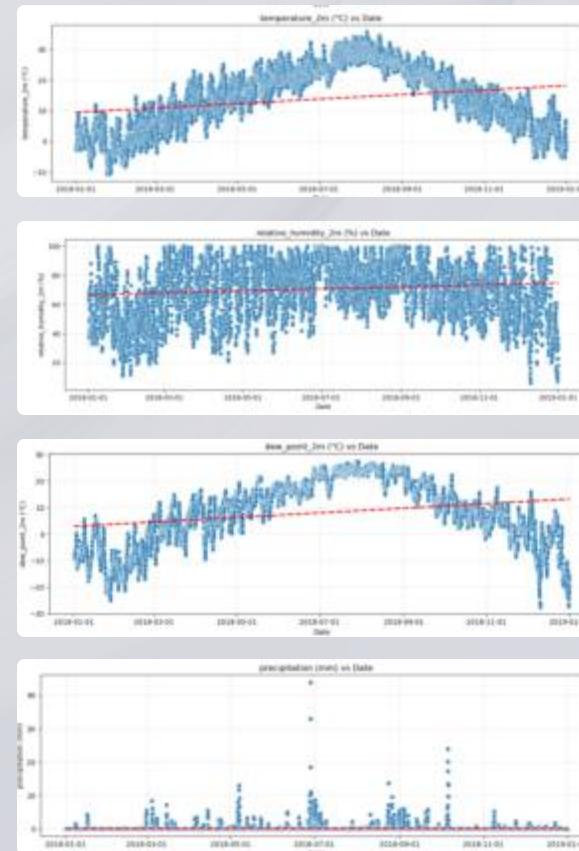
# TIME SERIES OF VARIOUS VARIABLE

.....and having a little fun with Gamma's AI animation

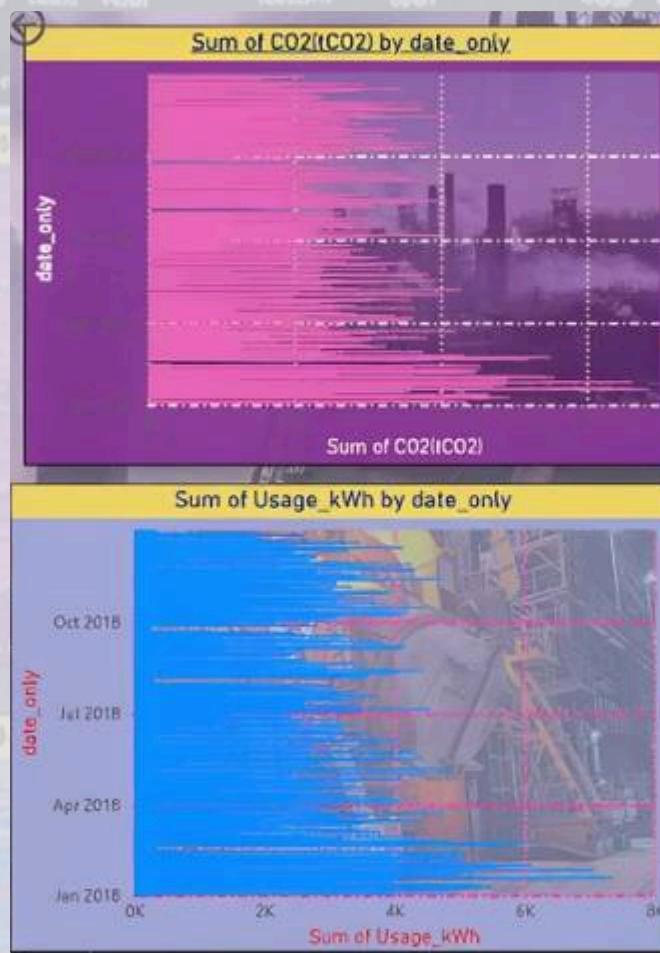
## Energy and CO2 variables



## Environmental /Weather variables

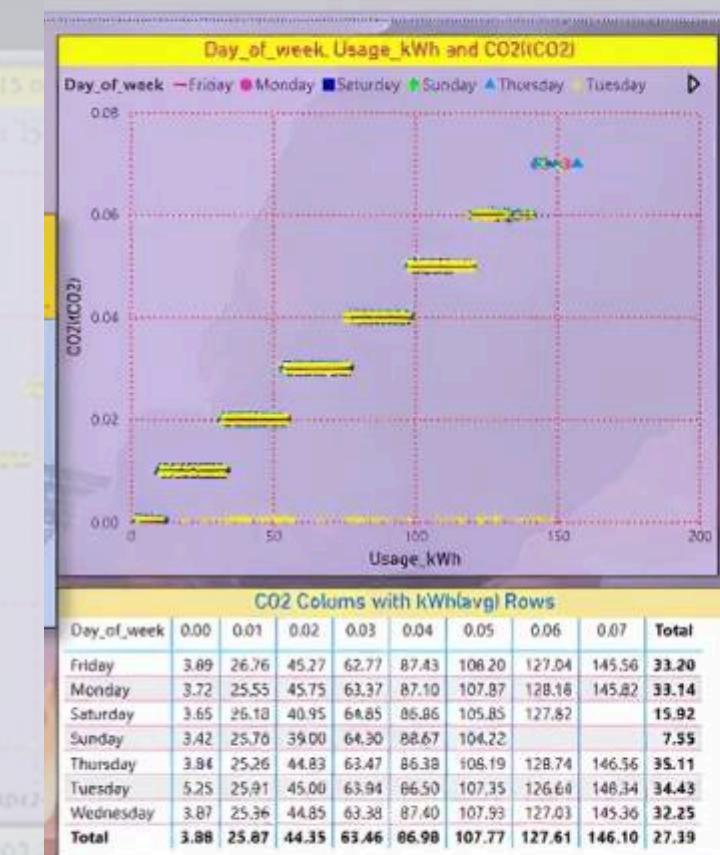


# Time Series Comparison of CO2 and Energy Usage kWh



## Yearly Look

The Graphs look almost identical with the trends. Both the CO2 emissions and Energy Usage kWh show peaks at the beginning winter months..

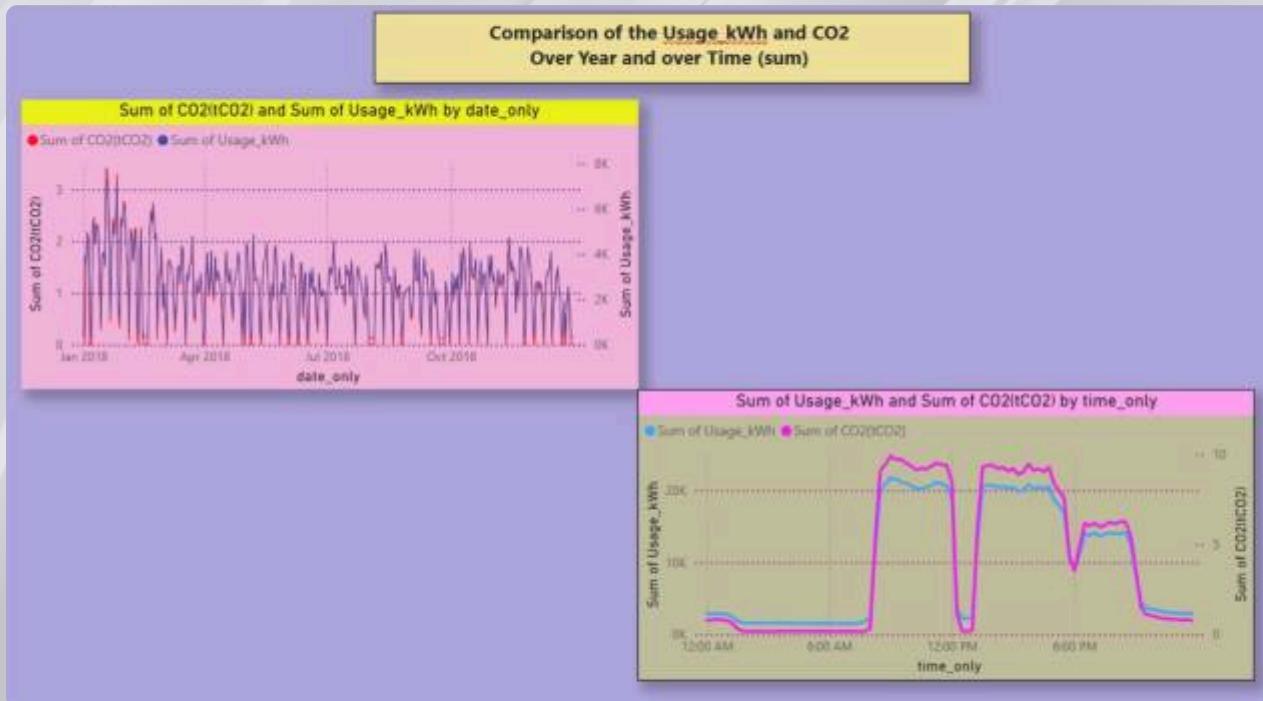


## Days of The Week

Scatterplot of CO2 vs. energy usage revealed a concentration of Tuesday observations (in yellow), suggesting a data volume skew – with Tuesdays contributing heavily to typical output ranges, but not to the highest-emission cases.”

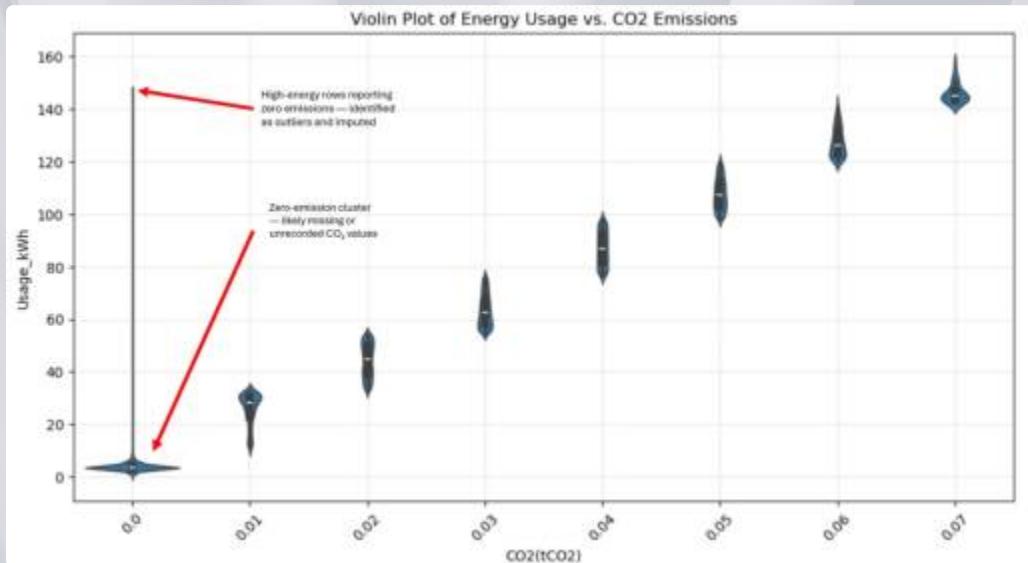
This correlation aligns with thermodynamic expectations – greater energy conversion yields higher CO<sub>2</sub> output. Also, it shows a large amount of zeros for the CO<sub>2</sub>. For some models Zero-emmission clusters removed during preprocessing and imputed for modeling consistency.

# Another Visual Comparing CO2 and Energy Usage kWh over time.



# Violin Plot

## Illustrates the Evolving Distribution of CO<sub>2</sub> across Energy Usage Levels



- Initial Variable Distributions">Histograms revealed skewed CO<sub>2</sub> emissions and wide energy usage range. Guided transformation choices (e.g., log\_kWh) and flagged missing values for imputation.

The violin plot reveals how CO<sub>2</sub> emission distributions change at different energy consumption levels, showing both frequency and probability density.

The plot revealed a dense spike at zero CO<sub>2</sub> emissions with a surprising vertical tail stretching into high energy usage levels. These outliers indicated likely data entry gaps since high kWh should yield measurable emissions. This informed targeted imputation to restore data integrity ahead of modeling.

The top energy level has wider areas where data points cluster most densely, indicating common operational states in the plant.

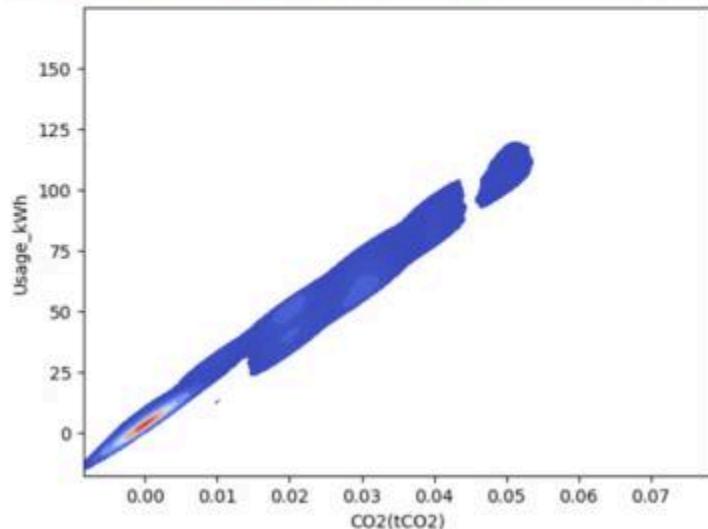
```
# Load necessary libraries
library(tidyverse)
library(magrittr)
library(vioplot)

# Load the data
air_filtration %>%
  na.omit() %>%
  mutate(Usage_kWh_imputed = ifelse(is.na(Usage_kWh), 0, Usage_kWh)) %>%
  ggplot(aes(x = CO2_tCO2, y = Usage_kWh_imputed)) +
  geom_violin() +
  geom_boxplot() +
  geom_jitter() +
  theme_minimal()

# Save the plot
ggsave("Violin_Plot_Energy_Usage_vs_CO2_imputed.pdf", width=10, height=6)
```

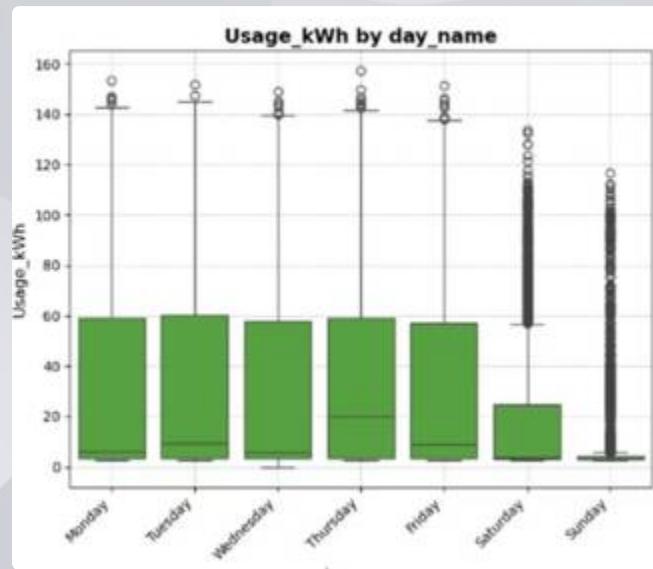
# Looking at the Density (Kernel)Plot to look

```
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
import numpy as np
#Density plot, compare the distribution of both variables
#had one same plot as the scatter and I didn't like it.
#Also, Updated syntax for newer versions of seaborn - using named parameters x and y instead of positional arguments
#And fill instead of shade for color too.
sns.kdeplot(x=df["CO2(tCO2)"], y=df["Usage_kWh"], cmap="coolwarm", fill=True)
plt.savefig("density_Plot_usage_KWh-CO2.jpg", dpi=300, bbox_inches="tight")
```



- The predominantly blue color with a small amount of red around intercept CO<sub>2</sub>, suggests that the highest density (red) is concentrated near the lower end of the CO<sub>2</sub> axis, with density decreasing (shifting to blue) as CO<sub>2</sub> increases.
- The thicker density from 0.02 to 0.05 CO<sub>2</sub>, followed by a split with no color, and then a thicker blob again at 0.06, indicates varying concentrations of data points. The "split with no color" likely represents an area of very low or zero density (a gap in the data distribution).
- The thicker portion of line indicates a higher concentration of data points, meaning more observations fall within that range of CO<sub>2</sub> (x-axis) and energy usage (y-axis).

# Boxplots illustrating how Energy Usage Fluctuates across time



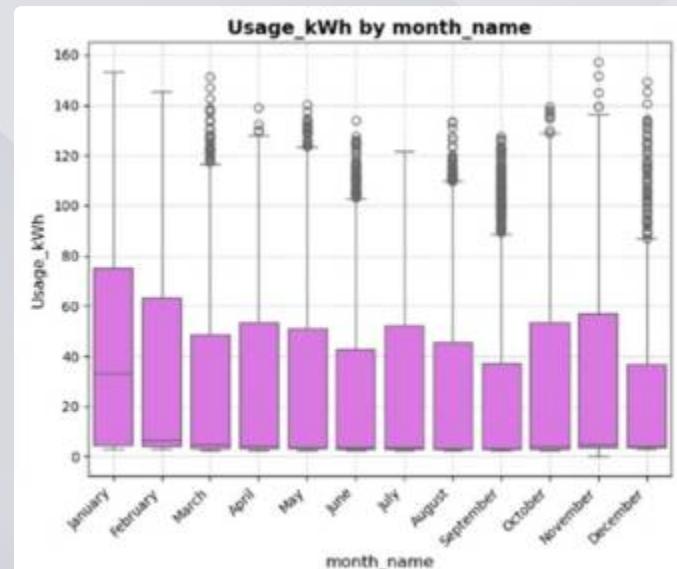
📅 Day of the Week

## Weekdays

No strong correlation across weekdays.

## Weekends

Shows **lower usage on weekends**, which is a clear and explainable pattern (possibly reduced shifts or partial operations/production activity). Later analysis, shows Tuesday as a possible variable to examine further.



📅 17 Monthly Patterns

**January & February stand out** as peak usage and aligns with heating needs or seasonal production cycles *"Winter peak usage reflects seasonal energy demand"*

**November shows early uptick**, possibly signaling seasonal prep or ramp-up

**December not matching upward trend** → plausible cause: holiday slowdowns or incomplete or partial data collection.

# Trend Analysis Over Time

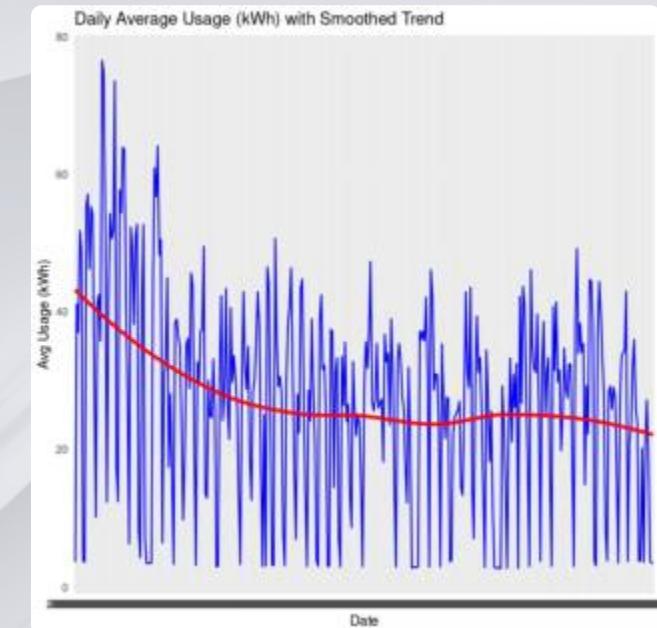
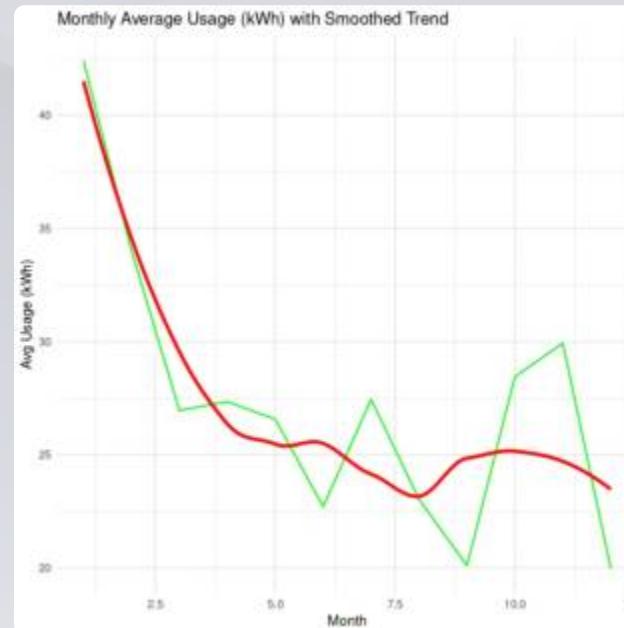
Supplements to the Boxplots showing monthly trends

Energy usage showed clear seasonal behavior, with the highest consumption occurring in the first two winter months (January and February). Usage declined steadily through spring and summer, followed by a noticeable spike in November. Interestingly, December did not continue this upward trend. It is likely a result of reduced plant activity during the holiday season or potential gaps in data collection.

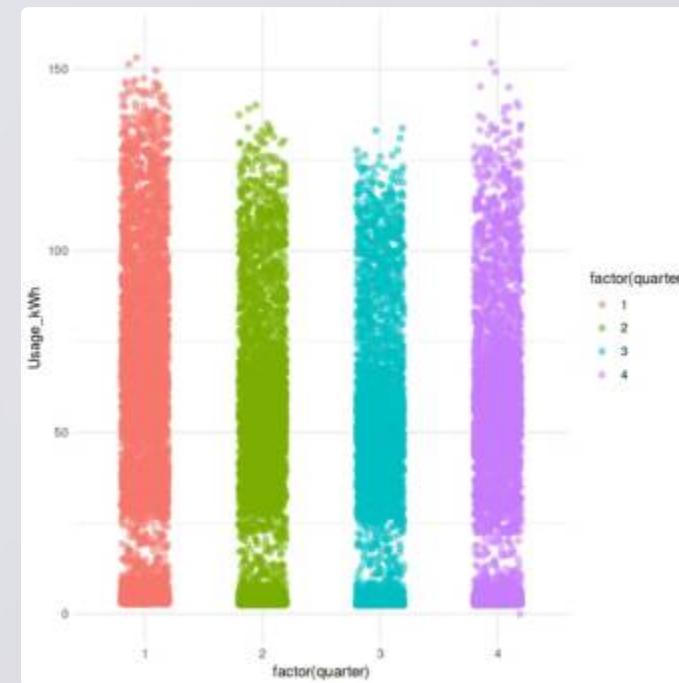
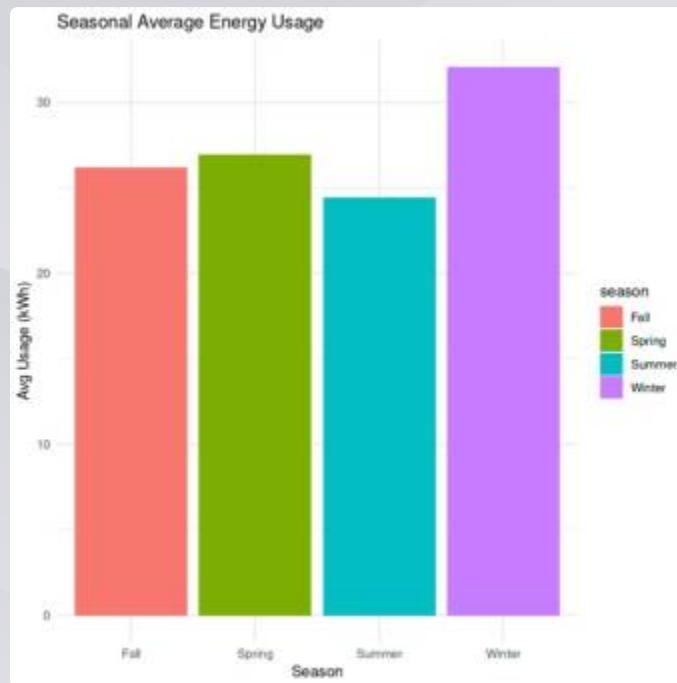
## CODE USING R IN KAGGLE

```
Looking at the monthly trend as well as time to plot monthly_avg_usage %>%  
  geom_line() + geom_smooth()  
  + theme_minimal() + labs( title = "Monthly Average Usage (kWh) with Smoothed Trend")
```

```
ggplot( monthly_avg_usage, aes( x = Month ) ) +  
  geom_line() + geom_smooth()  
  + theme_minimal() +  
  labs( title = "Monthly Average Usage (kWh) with Smoothed Trend")
```



# SEASONS and Energy Usage (kWh)



```
Play! Look at seasonal trends, not date if well above a trend, and fun to try library! | ← 100%  
mutate( | + case_when( | Month <= 12, 1, 2) ~ "Winter", Month >= 1, 4, 5) ~ "Spring", |  
Month <= 9, 10, 11) ~ "Fall", 0) | Month group_by( |) | 100%  
summarise( | = mean( |) | → 100%)
```

Consistent with other temporal patterns, the winter months show the highest average energy usage — reinforcing the seasonal impact on operational demand.

```
ggplot( |, aes( | = factor( |), | = factor( |, |) | = factor( |, |))) |  
geom_jitter( | = 0.2, | = 0, | = 0.5) | + theme_minimal()
```

Geo jitter plot by quarter illustrates seasonal energy use variability. Q1 shows the highest median usage and widest spread — consistent with winter operations and fluctuating demand. Q4 also displays variability, driven by a November spike followed by a December drop, likely due to holiday slowdowns or partial data capture.

# Anova with Energy Usage kWh with Environmental weather variables

As I began modeling, my exploratory analysis continued in parallel. While initial insights shaped early model choices, I remained open to discovering new patterns and refining features throughout the process. This iterative approach allowed for deeper understanding of environmental drivers and ensured each model was informed by the evolving data story.

## ANOVA For Energy Usage kWh with CO2

	sum_sq	df	F	PR(>F)
temperature_2m	5.689358e+02	1.0	2.145128e+01	3.042525e-06
relative_humidity_2m	2.773542e+03	1.0	3.866614e+02	7.746798e-25
precipitation	3.785748e+01	1.0	1.427717e+00	2.208991e-01
CO2	5.158338e+02	1.0	1.345873e+00	0.000000e+00
Residual	9.183427e+05	35035.0	Null	Null

## ANOVA For Energy Usage kWh without CO2

	sum_sq	df	F	PR(>F)
temperature_2m	6.268148e+02	1.0	400.194908	1.782512e-131
relative_humidity_2m	3.003326e+06	1.0	2911.016511	0.000000e+00
precipitation	6.217191e+04	1.0	40.917008	1.687528e-01
Residual	3.610988e+07	35036.0	Null	Null

Variance Inflation Factors:		
	Variable	VIF
0	const	16.677615
1	temperature_2m	1.153531
2	relative_humidity_2m	1.283317
3	precipitation	1.061145
4	CO2	1.082835

## Key Statistical Findings

### When CO<sub>2</sub> is included:

- CO<sub>2</sub> is dominant** in explaining energy usage ( $F = 1.35$  million,  $p < 0.00001$ ).
- Other variables like temperature and humidity still matter, but their influence pales in comparison.
- Precipitation becomes statistically irrelevant ( $p = 0.23$ ), likely overshadowed by the strength of CO<sub>2</sub>.

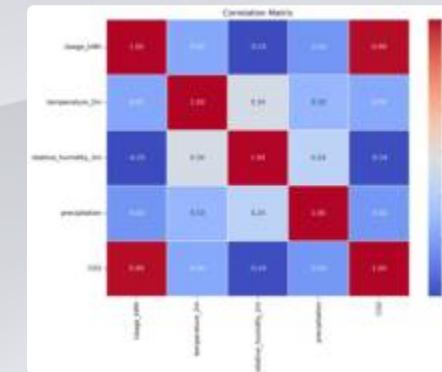
### When CO<sub>2</sub> is excluded:

- Relative humidity** becomes the strongest predictor ( $F = 2914$ ).
- Temperature** also has a sizable effect ( $F = 400$ ).
- Precipitation** is significant depending on the  $p$ -value formatting (though 1.6076 looks like a formatting error — likely  $p < 0.0001$  based on the F-stat).

## Interpretation:

CO<sub>2</sub> captures a large portion of the variance explained by both humidity and temperature. When removed, their independent effects surface more strongly.

VIF results confirm feature stability, with no signs of problematic multicollinearity.



**HEATMAP** of Energy Usage and CO<sub>2</sub> with Weather-related environmental variables — including temperature, humidity, and precipitation.

While the correlation between CO<sub>2</sub> and energy usage aligned with expectations, exploratory analysis revealed an inverse relationship between humidity and both target variables. This pattern suggests drier atmospheric conditions may correspond with higher energy demand and emissions. Temperature and humidity also exhibit interdependence, warranting deeper investigation through interaction modeling in subsequent analysis.

# Anova with CO2 with Environmental weather variables

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Define the formula with CO2 as the dependent variable
formula = "CO2 ~ Usage_kWh + temperature_2m + relative_humidity_2m + precipitation"

# Fit the ANOVA model
model = smf.ols(formula, data=df).fit()
anova_table = sm.stats.anova_lm(model, type=2) # Type 2 for balanced designs
print(anova_table)
```

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Define the formula with CO2 as the dependent variable
formula = "CO2 ~ temperature_2m + relative_humidity_2m + precipitation"

# Fit the ANOVA model
model = smf.ols(formula, data=df).fit()
anova_table = sm.stats.anova_lm(model, type=2) # Type 2 for balanced designs
print(anova_table)
```

## ANOVA with Energy Usage kWh

	sum_sq	df	F	PR(>F)
Usage_kWh	8.226553e+00	1.0	1.345871e+06	0.000000e+00
temperature_2m	4.562720e-04	1.0	7.464646e+01	5.864480e-18
relative_humidity_2m	2.564990e-05	1.0	4.196345e+00	4.051864e-02
precipitation	2.144567e-07	1.0	3.508529e-02	8.514181e-01
Residual	2.141493e-01	35035.0	NaN	NaN

## ANOVA without Energy Usage kWh

	sum_sq	df	F	PR(>F)
temperature_2m	0.159554	1.0	662.281719	1.058104e-144
relative_humidity_2m	0.675527	1.0	2804.004621	0.000000e+00
precipitation	0.009517	1.0	39.503582	3.313068e-10
Residual	8.440702	35036.0	NaN	NaN

## CO<sub>2</sub> as a Response Variable

- **Usage\_kWh is once again the strongest driver of CO<sub>2</sub> emissions ( $F \approx 1.35$  million,  $p < 0.00001$ ).**
- **Temperature & Humidity** also contribute, though with much smaller effect sizes.
- Precipitation remained the least influential among the weather variables

# Possible Interactions for CO2 with Energy Usage and Weather Variables

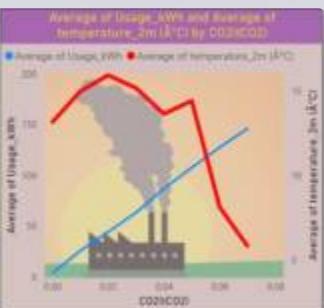
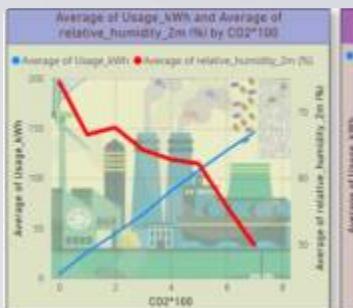
## ANOVA CHECK

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Define the formula with interaction terms
formula = "CO2 ~ usage_kWh + temperature_2m + usage_kWh * relative_humidity_2m + usage_kWh * precipitation"

# Fit the ANOVA model
model = smf.oaxa(formula, dataset).fit()
anova_table = sm.stats.anova_lm(model, typ=2)

print(anova_table)
```



	sum_sq	df	F \
Usage_kWh	8.226553e+00	1.0	1.348491e+06
temperature_2m	4.476525e-04	1.0	7.337892e+01
Usage_kWh:temperature_2m	3.474503e-05	1.0	5.695383e+00
relative_humidity_2m	2.218008e-05	1.0	3.635745e+00
Usage_kWh:relative_humidity_2m	4.227809e-04	1.0	6.930199e+01
precipitation	2.653875e-07	1.0	4.350216e-02
Usage_kWh:precipitation	1.543927e-05	1.0	2.530795e+00
Residual	2.137148e-01	35032.0	NaN

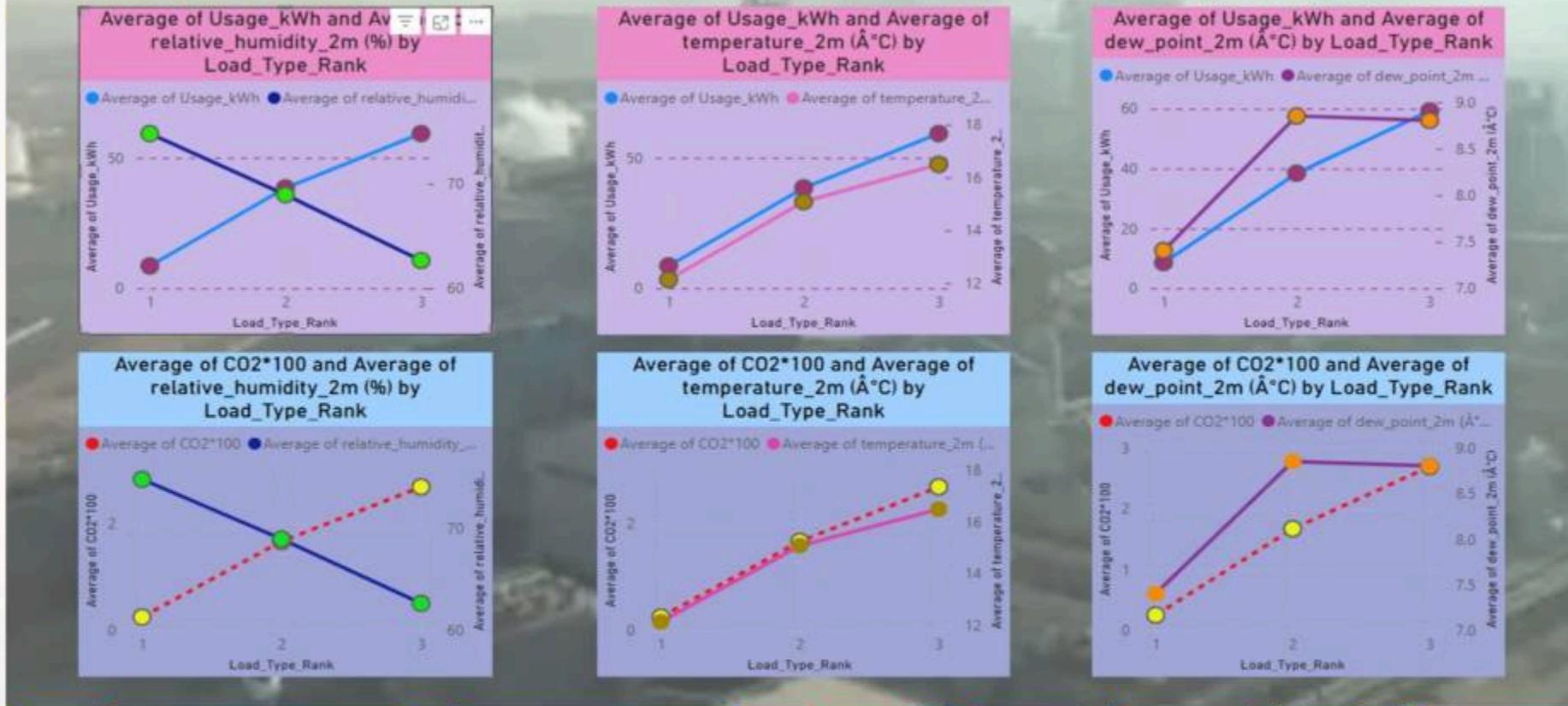
	PR(>F)
Usage_kWh	0.000000e+00
temperature_2m	1.113042e-17
Usage_kWh:temperature_2m	1.701487e-02
relative_humidity_2m	5.655948e-02
Usage_kWh:relative_humidity_2m	8.751208e-17
precipitation	8.347838e-01
Usage_kWh:precipitation	1.116530e-01
Residual	NaN

## 🔍 Interpretation Notes

- No Surprise, **Usage\_kWh is the dominant driver** of CO<sub>2</sub> variation
- **Interactions with temperature and humidity are crucial:** they reveal that energy usage affects CO<sub>2</sub> differently depending on these environmental conditions.
- **Precipitation seems negligible here**, both on its own and in interaction with Usage\_kWh.

## **Section 3: Feature Engineering and Diagnostics**

## Weather with CO2 and kWh Comparison



### Key Points:

- Energy and CO<sub>2</sub> move in lockstep, confirming physical expectations and past analysis
- Humidity's inverse trajectory highlights a **possible interaction effect**, where drier conditions align with higher energy and emissions with the crossing point at Rank 2 marking a shift in environmental influence.
- Dew point show a slight decoupled patterns, reinforcing that **drier air correlates with higher energy demand**
- Load Rank 2 consistently emerges as a pivot point, visually and statistically, where operational and environmental forces intersect

# HEATMAPS

## Humidity x Temp Interaction (CO2 and Usage kWh)

Humidity Bins by Temperature Bins											
humiditybinby10	-15.00	-10.00	-5.00	0.00	5.00	10.00	20.00	25.00	30.00	Total	
Total	1.5237	1.0548	1.0970	1.0906	1.1431	1.0721	1.1290	1.2871	1.8530	1.1524	
0			0.0000	0.5000							0.4286
10		0.3333	0.8616	1.6250							1.0357
20	2.7500	1.8147	1.3077	2.3530		2.8125	2.1667	2.6250			1.9038
30	2.7813	1.2270	1.8452	1.8806	2.7034	1.9750	2.9211	2.5714			1.9403
40	1.6445	1.4588	1.8406	1.3413	2.1763	2.0131	1.9511	1.2823	0.3438		1.7433
50	1.4670	1.0353	1.0193	1.6346	1.7172	2.0519	1.4733	1.8488	1.7599		1.5242
60	0.2541	0.5678	0.9167	1.2609	1.1852	1.2630	1.5654	1.8208	1.8581		1.2916
70		0.0000	0.3708	0.7128	0.8193	0.8733	0.9317	1.3678	1.6445	2.3638	1.0732
80		0.2100	0.9311	0.3827	0.7325	0.5851	1.0577	0.9168	2.6250		0.7443
90			0.8153	1.0096	0.6266	0.5688	0.7148	0.8796			0.7278
100			1.0000	0.0000	1.0769	0.2594	0.6637				0.6054

Humidity Bins by Temperature Bins												
humiditybinby10	-15.00	-10.00	-5.00	0.00	5.00	10.00	20.00	25.00	30.00	Total		
Total	34.41	25.56	26.73	26.77	27.10	25.54	26.48	29.81	41.81	27.39		
0			3.63	13.82							12.36	
10			13.85	21.39	37.81						25.42	
20			41.63	31.11	52.78		42.80	47.83	58.44		43.40	
30			23.90	29.57	42.79	44.25	51.24	34.16	33.00	37.21	44.37	
40			35.72	33.44	42.80	33.44	48.87	43.64	45.19	30.37	11.58	40.20
50			34.06	25.10	25.05	34.62	39.51	46.06	33.29	42.06	39.93	35.10
60			10.72	13.18	22.53	30.54	26.36	29.37	35.43	40.26	41.71	30.22
70			4.89	16.32	18.93	20.42	21.47	22.35	32.09	36.99	33.01	25.66
80			8.43	23.47	11.57	18.54	15.39	24.84	22.55	34.45		18.80
90			20.41	24.81	16.12	15.12	17.98	21.28				18.34
100			36.07	3.36	25.90	8.65	16.84					15.80

### Analysis and Key Takeaways

#### Hot Zones of Inefficiency

Heatmaps reveal two distinct emission and energy usage spikes: one in moderate humidity (20–40%) and another in high humidity (70–80%) at elevated temperatures. These zones may reflect both optimal production and environmental strain which is highlighting the need for targeted ventilation or load-balancing strategies.

# Energy Variables and the Definitions

While the initial focus centered on environmental weather variables, the heatmap reveals strong correlations with several energy measurements from the steel plant. These patterns are too meaningful to ignore. It's time to explore those energy metrics in more depth and integrate them into the upcoming models.

## Lagging Current Reactive Power (kVarh)

Represents reactive power consumed by **inductive loads** (e.g., motors, transformers).

Does not contribute directly to useful work but impacts **system efficiency**.

Higher lagging reactive power leads to lower power factor and increased energy losses.

## Lagging Current Power Factor

Power factor is the ratio of **real power (kWh)** to **apparent power (kVA)**.

A **lagging power factor** occurs when current lags behind voltage due to inductive behavior.

Lower power factor means more wasted energy and higher electricity costs.

## Leading Current Reactive Power (kVarh)

Represents reactive power consumed by **capacitive loads** (e.g., capacitor banks).

Occurs when current leads voltage, opposite of inductive loads.

Excessive leading reactive power can also reduce system efficiency.

## Leading Current Power Factor

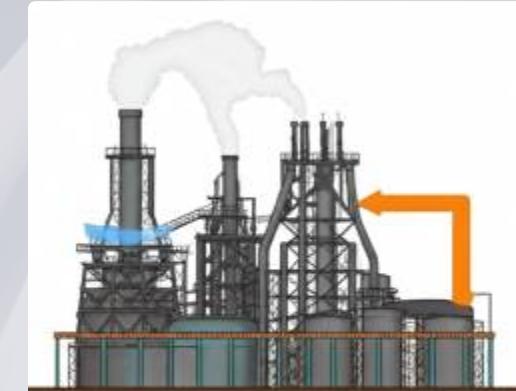
Measures the extent to which current leads voltage in a capacitive circuit.

Indicates **capacitive dominance** in the power system.

Excessively leading power factor may **overcorrect system behavior**, causing:

Voltage rise, Resonance issues, Reduced efficiency

May require reactive inductive loads (e.g., reactors) to stabilize the system.



## Energy Usage (kWh)

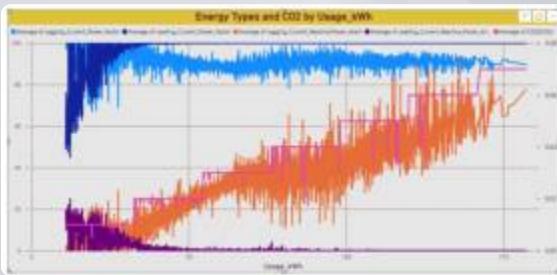
Represents the total amount of **real energy consumed** by the steel plant over time while capturing the core operational demand across all machinery and processes.

## Load Rank (1–3)

1	<b>Light Load</b> – Low production state, typically associated with minimal energy demand.
2	<b>Medium Load</b> – Standard production cycles; balanced usage across most equipment. Reflects routine activity levels.
3	<b>Max Load</b> – Peak energy demand; full-scale production with all systems engaged.

# Energy Variables Graphs

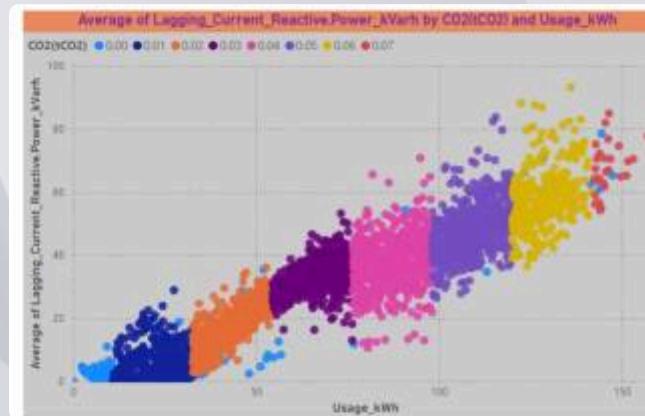
## Multi-Variable Energy Scatter: Mapping Power Behavior and CO<sub>2</sub> Emissions



The **linear structure** underscores a steady relationship between usage and reactive load (Lagging Reactive Power kVarh Variable), indicating inductive demand behaving predictably.

The CO<sub>2</sub> follows the same trend and is on top of the Lagging Reactive Power.

## Energy Usage, Reactive Load, and CO<sub>2</sub>



**Color gradations represent CO<sub>2</sub> intensity**, with warmer hues signaling rising emissions

The **stacked segments** hint at discrete operational modes and each with its own CO<sub>2</sub> footprint

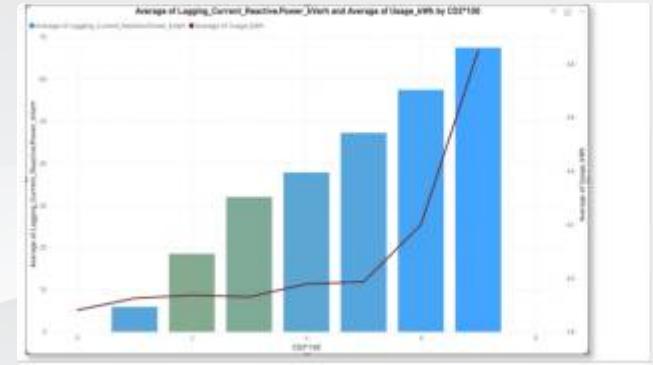
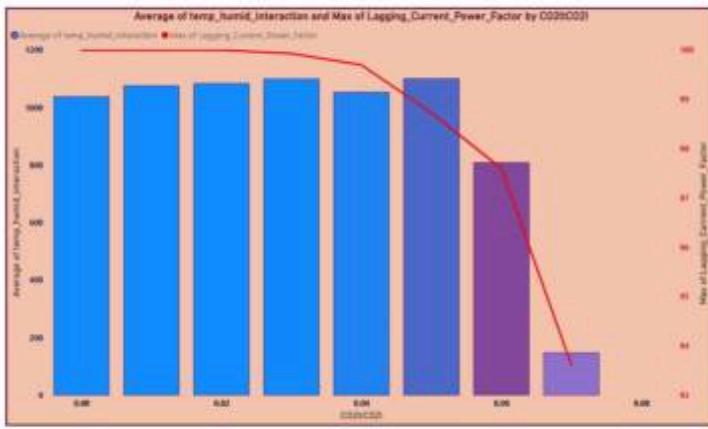
Outliers in light blue scattered throughout may be imputed gaps and unexplained behavior.

## The Visual of Steelplant Dynamics

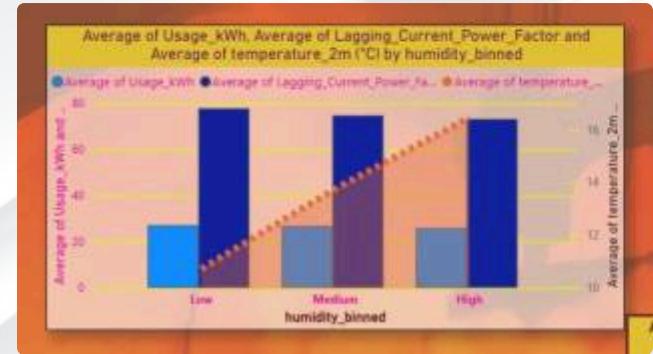
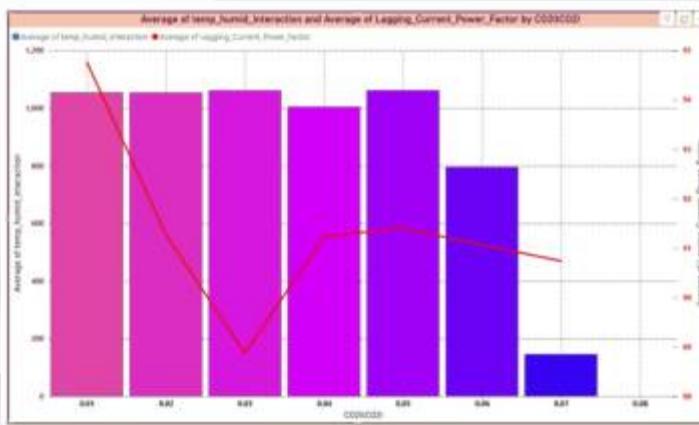
Each color stratum performs its part in the energy-emissions ensemble. The two Scatter plots charts **Usage kWh** against **Lagging Reactive Power kVarh**, forming a linear composition. The **color-coded bands**, segmented by CO<sub>2</sub> levels from light blue (0.0) to red (~0.07), reveal distinct operational layers across the plant's load profile.

# Lagging Current Reactive Power kVarh variable

## Energy Behaviors Leads Emission Shifts with environmental weather influence

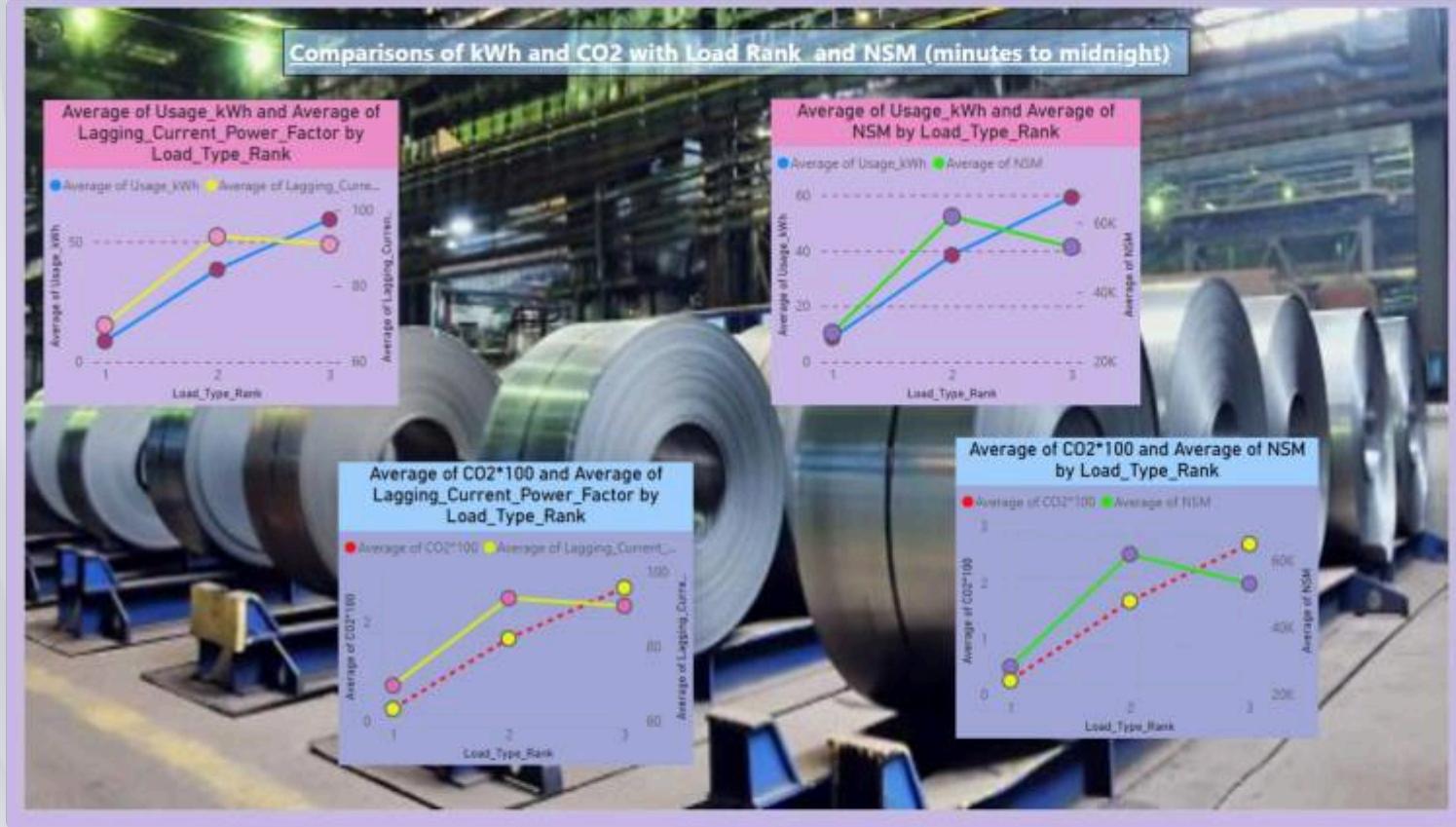


Both lagging energy and usage rise with CO<sub>2</sub>, but there's a *nonlinear surge* at the highest emission levels.



Energy usage falls as humidity (and temperature) rise and is possibly cooling demand. The Cooling Effect of Moisture: Reactive and Real Energy Dip with Humidity?

- The *max lagging* graph shows a **sharp drop** at high CO<sub>2</sub>, implying system breakdowns or constraint thresholds.
- The *average lagging* graph adds **nuance**: it reveals that the system *partially rebounds* but doesn't return to its initial state, suggesting **nonlinear or hysteretic response**.
- Together, they suggest there's both a *disruption point* and a *stabilized inefficiency*, which is gold in a data story.



## Key points

- While energy usage climbs steadily with load rank, the timing of those loads' shifts; medium loads peak later, while max loads occur slightly earlier in the day.
- The Lagging Current Power Factor increases at level 2 and then also levels out, again could be due to timing of load shifts.



## Key points

- The Lagging Current Reactive Power shows correlation system load, not as strong as Energy usage or CO<sub>2</sub>
- High CO<sub>2</sub> emission periods demonstrate sharp drops in lagging power, suggesting system constraints or operational thresholds.
- Medium loads demonstrate different timing patterns compared to maximum loads, affecting the overall power factor efficiency.
- The data reveals both disruption points and stabilized inefficiency regions, providing valuable insights for optimization strategies.



## Lagging Power Factor by Load Rank

Comparing two key electrical power metrics and their behavior patterns



### Lagging Reactive Power (kVArh)

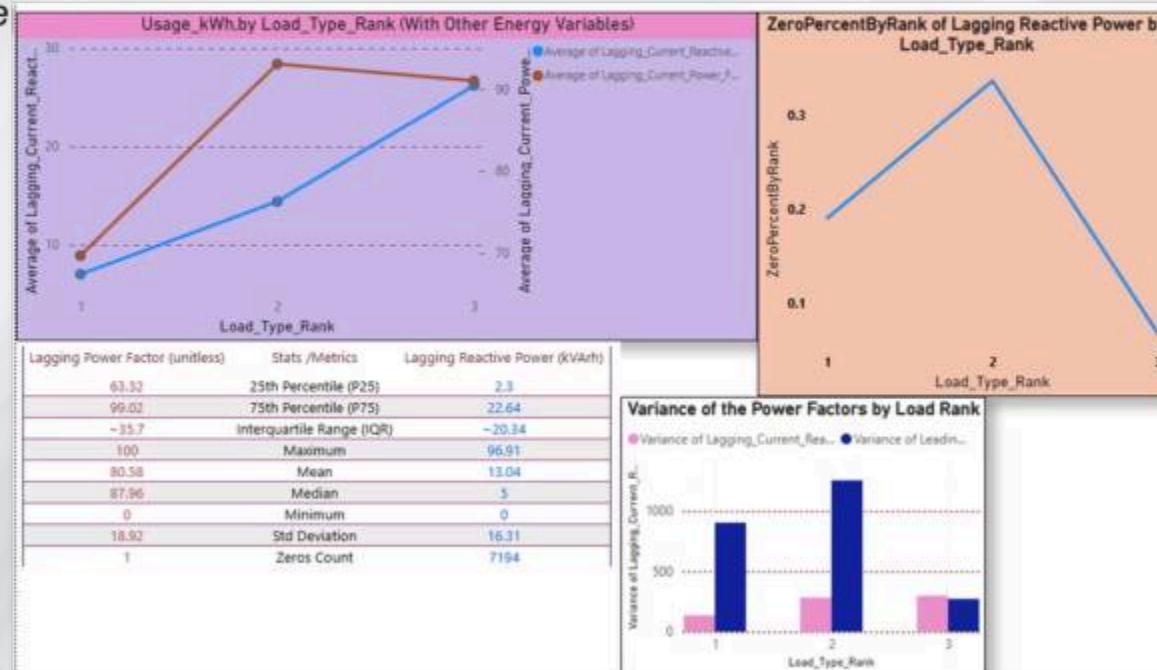
**Definition:** Total reactive energy used (inductive demand)

**Load Rank Trend:** Increases steadily (~45° linear climb)

**Zero Values:** High (7,194) zeros and many load periods with draw no reactive energy

**Distribution:** Right-skewed with P25 = 2.30, P75 = 22.64

**Variability:** Standard deviation ~16.3 (moderate spread)



### Lagging Power Factor (unitless)

**Definition:** Efficiency of power usage ( $\text{kWh} \div \text{kVA}$ )

**Load Rank Trend:** Starts near zero, rises quickly, levels by rank 3

**Zero Values:** Only 1 zero and PF typically registers even at low loads

**Distribution:** Wide efficiency range with P25 = 63.32, P75 = 99.02

**Variability:** Standard deviation ~18.9 (broad range; nonlinear distribution)

## Section 4: Modeling Landscape: Structural, Temporal, and Behavioral Insights

*While not every model is featured here, each played a role in shaping the final lineup. The modeling journey was iterative, exploratory, and occasionally messy, but it led to the clarity and performance showcased in the next slides. They are a distilled selection from a much larger pool of models and diagnostics still included in the portfolio.*

*Early SARIMAX and ARIMA fits guided much of the rhythm-based exploration, even though they aren't fully presented due to memory limitations and partial diagnostics. Their influence is woven into the time-aware analysis that follow and are reflected in weekday/weekend behaviors, temporal features, and log\_kWh modeling.*

# Model REady?

## OLS Regression Model for CO<sub>2</sub>

As I began modeling, my exploratory analysis continued in parallel. While initial insights shaped early model choices, I remained open to discovering new patterns and refining features throughout the process. This iterative approach allowed for deeper understanding of environmental drivers and ensured each model was informed by the evolving data story.

First model with focus on the Environmental Weather Variables and Focus

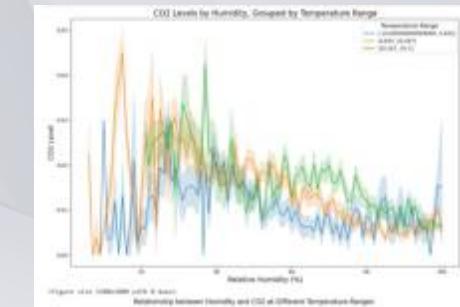
OLS Regression Results					
Model Variable	Estimate	Standard Error	t Statistic	p-value	AIC
Model:	Least Squares	F-statistic(1, 1)	4.375e+00	0.00	
Date:	04/18/2023	F-statistic(1, 1)	4.375e+00	0.00	
Days:	33-2773	Log Likelihood	-1.321e+00		
N Observations:	33668	AIC	-2.643e+00		
DF Residuals:	33667	BIC	-1.443e+00		
DF Model:	1				
Convergence Type:	converged				
Coefficients:					
Intercept	0.0000	8.37e-01	-0.235	0.81%	
Temperature_2m_degC	0.0000	8.37e-01	-0.235	0.81%	
Logging_Current_Reactive_Power_kVarh	0.0000	8.37e-01	-0.235	0.81%	
Temp_Humid_Interaction	-0.423e-06	2.39e-07	18.365	0.000	-1.1e-06
Standard Errors:					
Deg2Kst	0.001100	0.000000	-0.407		
StdErr(Beta0):	0.000400	0.000000	-0.000		29984.140
StdErr(Beta1):	0.000000	0.000000	-0.000		0.000
StdErr(Beta2):	0.000000	0.000000	-0.000		0.000
Residuals:					
Residuals:	0.000000	0.000000	0.000000	0.000000	0.000000
Notes:					
231 Standard Errors assume that the covariance matrix of the errors is correctly specified.					
Interaction term p-value: 4.48837712988084e-05					

### 🔍 Interpretation Notes

- Model explains ~78.9% of variance in CO<sub>2</sub>
- The interaction effect is powerfully negative, indicating that **high humidity dampens the temperature-driven CO<sub>2</sub> rise**, possibly due to altered ventilation or energy dynamics.
- The t-statistic for reactive power is high and it's an important driver.
- Durbin-Watson (~0.41)** shows strong positive autocorrelation in residuals and maybe CO<sub>2</sub> has time-lagged patterns. Will address via ARIMA/SARIMA models later.

### Regression Highlights

Variable	Coef	P-value	Insight
Intercept	-0.0004	< 0.001	Slight negative baseline CO <sub>2</sub> offset
Temp (°C)	3.119e-05	< 0.001	↑ Temp → slight increase in CO <sub>2</sub>
Reactive Power (kVarh)	0.0009	≈ 0	Strong positive CO <sub>2</sub> association
Temp × Humidity Interaction	-4.63e-06	≈ 0	Highly significant suppressive effect



# OLS Regression Model Continues for CO<sub>2</sub>

Continue to understand environmental variables and how much additional variance is explained by the power consumption variables.

```
# Step 2: Add power variables one by one, keeping those that improve the model
power_vars = [
    'Lagging_Current_Reactive_Power_kVarh',
    'Leading_Current_Reactive_Power_kVarh',
    'Lagging_Current_Power_Factor',
    'Leading_Current_Power_Factor']
```

Final Model Summary:

## OLS Regression Results

```
=====
Dep. Variable: CO2(tCO2) R-squared:      0.897
Model:          OLS   Adj. R-squared:  0.897
Method:         Least Squares F-statistic:   5.106e+04
Date: Fri, 20 Jun 2025 Prob (F-statistic): 0.00
Time: 19:27:53 Log-Likelihood: 1.3474e+05
No. Observations: 35040 AIC:        -2.695e+05
Df Residuals: 35033 BIC:        -2.694e+05
Df Model: 6
Covariance Type: nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	-0.0381	0.000	-127.056	0.000	-0.039	-0.038
temperature_2m (°C)	-0.0001	1.06e-05	-9.898	0.000	-0.000	-8.41e-05
relative_humidity_2m (%)	-5.425e-05	2.44e-06	-22.247	0.000	-5.9e-05	-4.95e-05
temp_humid_interaction	2.486e-06	1.53e-07	16.295	0.000	2.19e-06	2.78e-06
Lagging_Current_Reactive.Power_kVarh	0.0007	2.13e-06	322.749	0.000	0.001	0.001
Lagging_Current_Power_Factor	0.0004	2.05e-06	179.583	0.000	0.000	0.000
Leading_Current_Power_Factor	0.0002	1.32e-06	123.627	0.000	0.000	0.000
Omnibus:	4711.603	Durbin-Watson:	0.520			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	61639.818			
Skew:	0.103	Prob(JB):	0.00			
Kurtosis:	9.494	Cond. No.	1.44e+04			

### Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.44e+04. This might indicate that there are strong multicollinearity or other numerical problems.

Will multicollinearity be reduced with added different Energy variables?

### Variance Inflation Factors for Final Model:

	Variable	VIF
0	const	117.730342
1	temperature_2m (°C)	14.726059
2	relative_humidity_2m (%)	3.099169
3	temp_humid_interaction	19.660000
4	Lagging_Current_Reactive.Power_kVarh	1.577807
5	Lagging_Current_Power_Factor	1.959951
6	Leading_Current_Power_Factor	2.112516

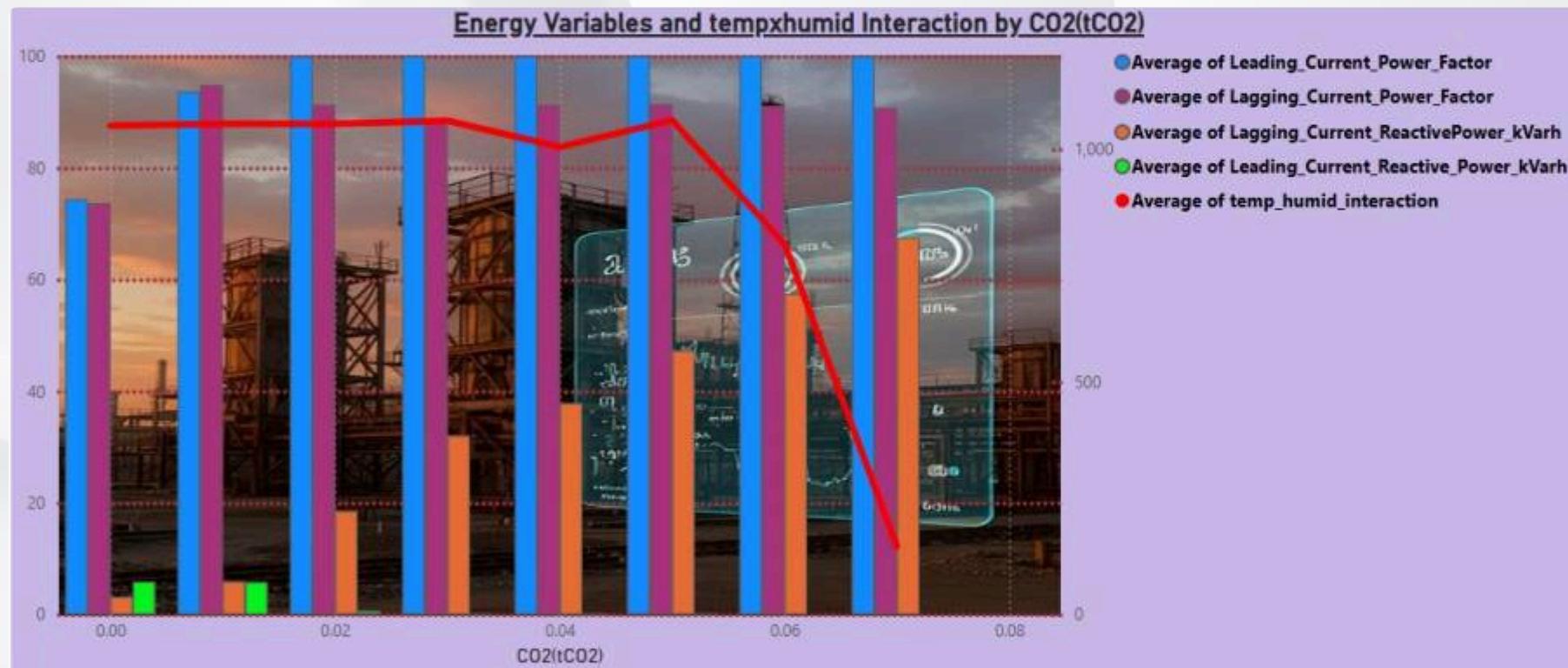
## Interpretation

- The model explains ~90% of the variance in CO<sub>2</sub>, driven primarily by power factors and reactive draw.
- While some multicollinearity exists among weather terms, especially the tempxhumidity interaction, it contributes significant explanatory power.
- Durbin-Watson value (~0.52) and non-normality hints (JB test) suggest residual autocorrelation and kurtosis
- Diagnostic indicators such as VIF remain within tolerable bounds for most variables, suggesting model stability.

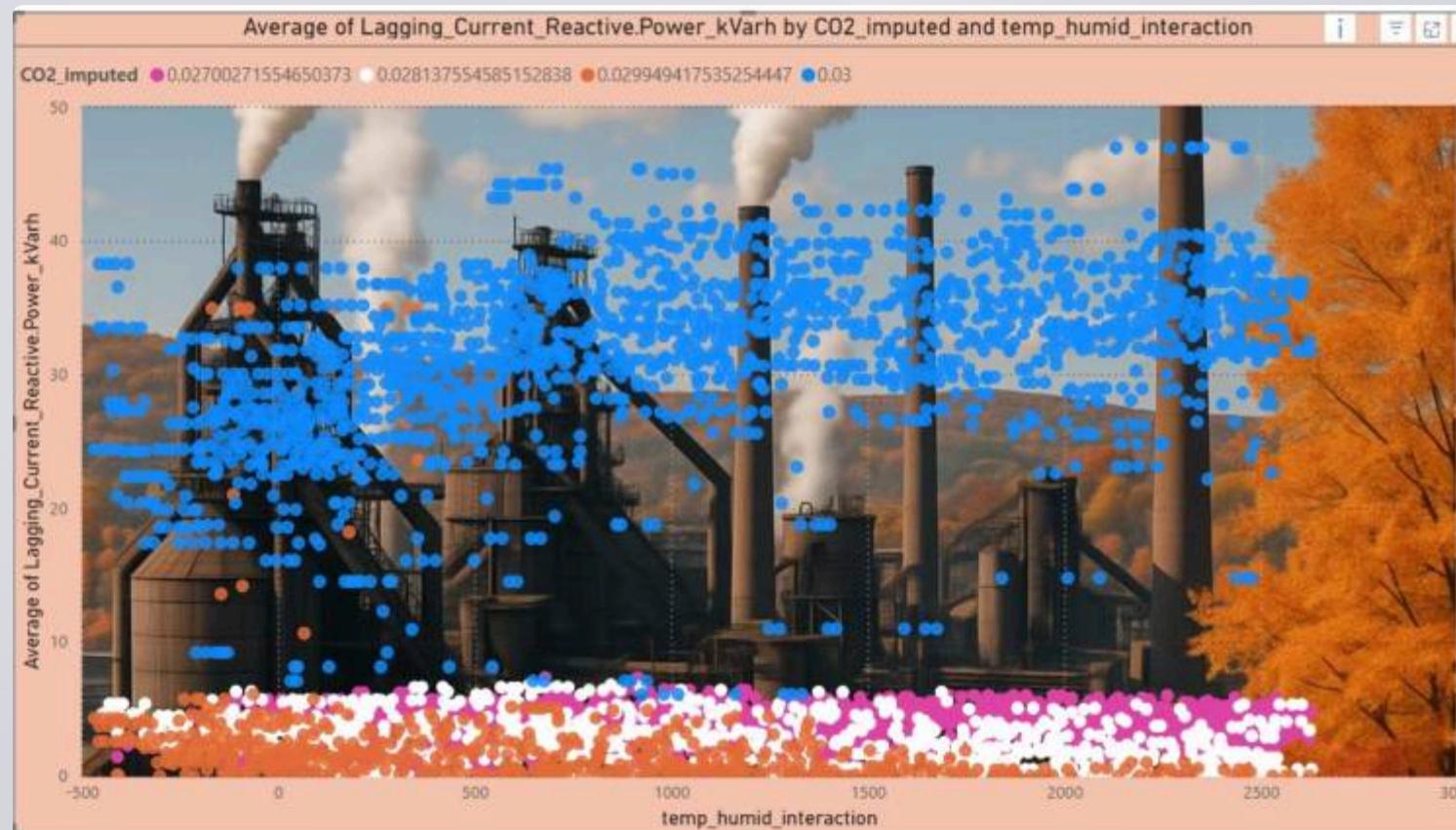
## Power vs. CO<sub>2</sub> Levels (TempxHumid)

**Observation:** While most power metrics increase predictably with rising CO<sub>2</sub> levels, leading current reactive power behaves differently. It remains suppressed at 0.00 and 0.01, barely emerges at 0.02, and then disappears again. This discontinuity may reflect system constraints, sensor thresholds, or maybe binning artifacts.

**Interpretation:** The faint emergence at 0.02 suggests a nonlinear threshold—possibly the point where capacitive behavior becomes detectable. Its absence at other levels may indicate operational suppression or measurement rounding.

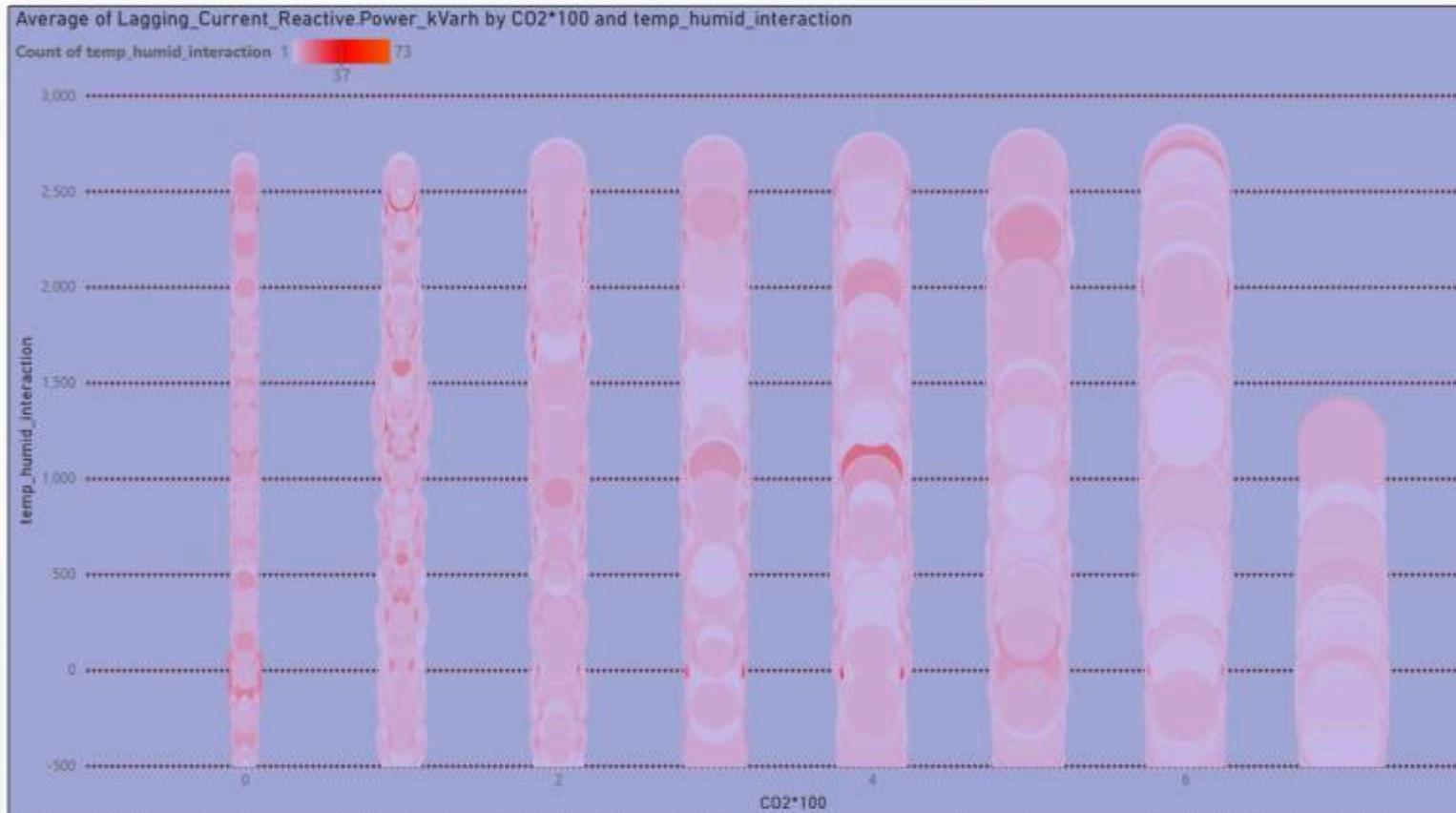


# Diagnostic Interpretation: Lagging Reactive Power vs. Imputed CO<sub>2</sub>



**Observation:** Across imputed CO<sub>2</sub> levels (0.00 to 0.03), lagging current reactive power remains low, except at **0.03**, where it spikes and clusters at higher values. This suggests a **threshold effect**, where the system shifts behavior once CO<sub>2</sub> reaches a critical level.

**Environmental Interaction:** Despite expectations, the temp × humidity interaction appears **non-influential** in this context. This reinforces the idea that **internal system load or control logic**, not environmental conditions, is driving reactive power behavior at high CO<sub>2</sub>.



## Stacked Signals: Reactive Power Density Across CO<sub>2</sub> Regimes

This stacked scatter reveals nonlinear shifts in reactive power behavior.

- Low CO<sub>2</sub> levels show minimal engagement
  - Mid-range levels (4–6) reveal sustained inefficiency zones
  - Level 7 suggests a short-duration spike worth diagnostic review
- Environmental interaction appears secondary and reactive behavior is primarily CO<sub>2</sub>-driven.

## OLS Regression Model Continues for Energy Usage kWh Following the same model as previous CO<sub>2</sub>

Environment Variables Model with kWh:

### OLS Regression Results

Dep. Variable:	Usage_kWh	R-squared:	0.078			
Model:	OLS	Adj. R-squared:	0.078			
Method:	Least Squares	F-statistic:	994.1			
Date:	Fri, 20 Jun 2025	Prob (F-statistic):	0.00			
Time:	19:27:53	Log-Likelihood:	-1.7127e+05			
No. Observations:	35040	AIC:	3.426e+05			
Df Residuals:	35036	BIC:	3.426e+05			
Df Model:	3					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	51.9243	0.941	55.165	0.000	50.079	53.769
temperature_2m (°C)	0.8115	0.065	12.462	0.000	0.684	0.939
relative_humidity_2m (%)	-0.4258	0.015	-28.846	0.000	-0.455	-0.397
temp_humid_interaction	-0.0054	0.001	-5.761	0.000	-0.007	-0.004
Omnibus:	5439.569	Durbin-Watson:	0.191			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	8319.668			
Skew:	1.151	Prob(JB):	0.00			
Kurtosis:	3.631	Cond. No.	7.26e+03			

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 7.26e+03. This might indicate that there are strong multicollinearity or other numerical problems.

Results of adding each power variable with y as kWh:

	Variable	AIC	R-squared	\
0	Lagging_Current_Reactive_Power_kVarh	285385.626440	0.819736	
1	Leading_Current_Reactive_Power_kVarh	285362.129974	0.819867	
2	Lagging_Current_Power_Factor	265106.554013	0.898955	
3	Leading_Current_Power_Factor	260115.366493	0.912375	

### Final Model Summary:

#### OLS Regression Results

Dep. Variable:	Usage_kWh	R-squared:	0.912			
Model:	OLS	Adj. R-squared:	0.912			
Method:	Least Squares	F-statistic:	5.211e+08			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	-1.3005e+05			
No. Observations:	35040	AIC:	2.601e+05			
Df Residuals:	35032	BIC:	2.602e+05			
Df Model:	7					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-80.5542	0.737	-109.365	0.000	-81.998	-79.111
temperature_2m (°C)	-0.2154	0.020	-10.628	0.000	-0.255	-0.176
relative_humidity_2m (%)	-0.1100	0.005	-23.576	0.000	-0.119	-0.101
temp_humid_interaction	0.0048	0.000	16.317	0.000	0.004	0.005
Lagging_Current_Reactive_Power_kVarh	1.4679	0.004	357.526	0.000	1.460	1.476
Leading_Current_Reactive_Power_kVarh	0.3667	0.022	16.556	0.000	0.323	0.410
Lagging_Current_Power_Factor	0.7395	0.004	186.807	0.000	0.732	0.747
Leading_Current_Power_Factor	0.3979	0.005	73.248	0.000	0.387	0.409
Omnibus:	5568.828	Durbin-Watson:	0.378			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	19907.101			
Skew:	0.784	Prob(JB):	0.00			
Kurtosis:	6.343	Cond. No.	1.85e+04			
	Variable	VIF				
0	const	193.918969				
1	temperature_2m (°C)	14.739592				
2	relative_humidity_2m (%)	3.099453				
3	temp_humid_interaction	19.660012				
4	Lagging_Current_Reactive_Power_kVarh	1.601896				
5	Leading_Current_Reactive_Power_kVarh	9.667420				
6	Lagging_Current_Power_Factor	2.005182				
7	Leading_Current_Power_Factor	9.785265				



**\*\*Interpretation of Model output on the next slide**

## OLS Regression for Energy Usage kWh

### Model Performance Interpretation



#### Environment Variables Only

R<sup>2</sup> Score: 0.078

AIC: 342,600

Very low explanatory power — weather alone doesn't account for energy demand shifts



#### Add Power Factors Incrementally

R<sup>2</sup> Score: Up to 0.912

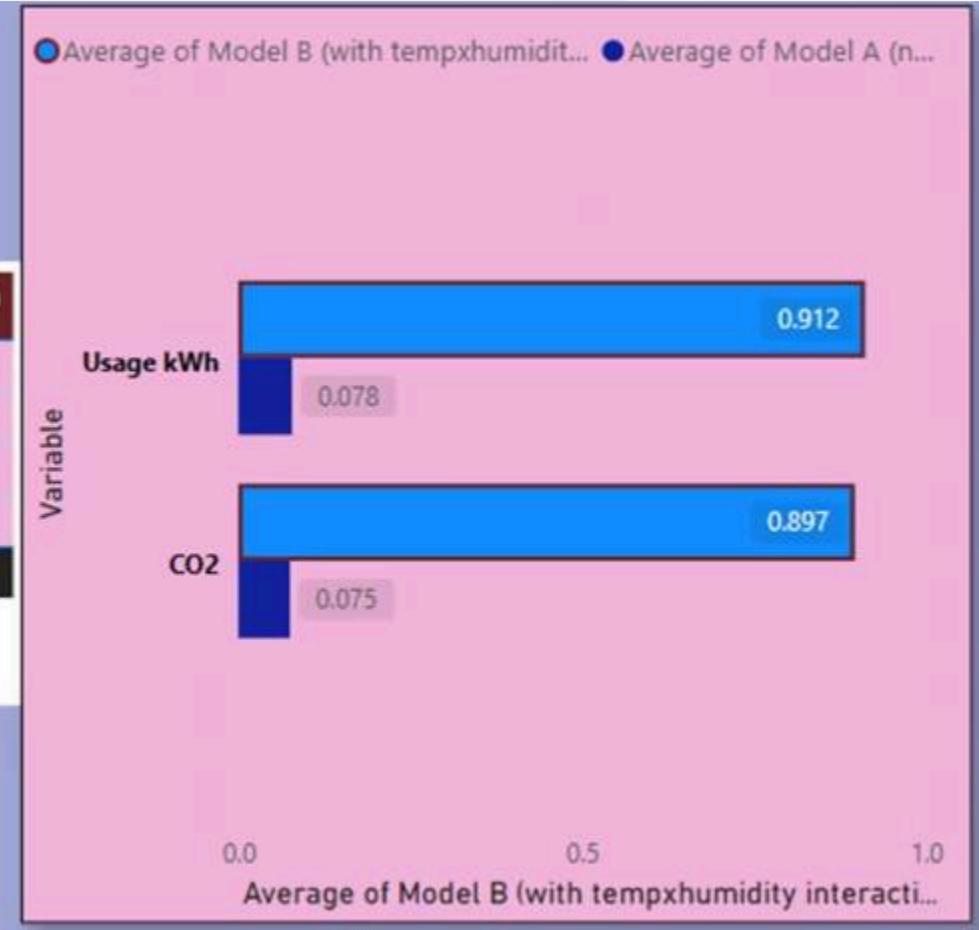
AIC: ↓ to 260,100

Dramatic boost — power dynamics dominate energy usage patterns

Adding power variables to the environmental baseline yields a 12x increase in explanatory strength (R<sup>2</sup> from 0.078 to 0.912). Among them, reactive power and power factor dominate, reinforcing the operational rather than environmental drivers of energy use. High VIFs among interaction and weather terms reflect their expected overlap, but don't compromise model validity.

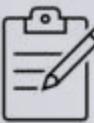
## Overview of with and without Interaction (tempxhumidity)

Metric	Model A (no interaction)	Model B (with Interaction)
R <sup>2</sup> (CO <sub>2</sub> )	0.075	0.897
R <sup>2</sup> (Usage_kWh)	0.078	0.912
AIC (CO <sub>2</sub> )	-192,439	-269,194
AIC (Usage_kWh)	342,587	260,115
Total		



# OLS Regression Model for CO<sub>2</sub>

## Dropping the tempxhumid interaction



NOTES: With the High VIF for tempxhumid interaction (19.66) and temperature\_2m (14.73) suggesting multicollinearity, I decided to drop the tempxhumid interaction, but will revisit later and maybe standardize the value. For the final model, I also dropped the Leading\_Current\_Reactive\_Power\_kVarh based on its lower coefficients and higher p-value.

### Environmental Variables Model (\*without interaction of tempxhumid\*):

#### OLS Regression Results

Dep. Variable:	CO2(tCO2)	R-squared:	0.075			
Model:	OLS	Adj. R-squared:	0.075			
Method:	Least Squares	F-statistic:	1430.			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	-96223.			
No. Observations:	35040	AIC:	-1.924e+05			
Df Residuals:	35037	BIC:	-1.924e+05			
Df Model:	3					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	0.0249	0.000	81.398	0.000	0.024	0.025
temperature_2m (°C)	0.0002	8.81e-06	25.838	0.000	0.000	0.000
relative_humidity_2m (%)	-0.0002	4.42e-06	-52.825	0.000	-0.000	-0.000
Omnibus:	4964.986	Durbin-Watson:	0.207			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	7345.697			
Skew:	1.099	Prob(JB):	0.00			
Kurtosis:	3.446	Cond. No.	276.			

#### Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

R-squared: 0.07545718474560481

AIC: -102439.8530152099

Power variables found in dataframe: ['Lagging\_Current\_Reactive\_Power\_kVarh', 'Leading\_Current\_Reactive\_Power\_kVarh', 'Lagging\_Current\_Power\_Factor', 'Leading\_Current\_Power\_Factor']

### Results of adding each power variable:

	Variable	AIC	R-squared
0	Lagging_Current_Reactive_Power_kVarh	-246448.152903	0.802071
1	Leading_Current_Reactive_Power_kVarh	-246446.611145	0.802073
2	Lagging_Current_Power_Factor	-256644.940135	0.852054
3	Leading_Current_Power_Factor	-269194.210347	0.896595

### Coefficient P-value

	Coefficient	P-value
0	0.0000859	0.00000
1	-0.000004	0.49848
2	0.000210	0.00000
3	0.000163	0.00000

### Final Model Summary:

#### OLS Regression Results

Dep. Variable:	CO2(tCO2)	R-squared:	0.897
Model:	OLS	Adj. R-squared:	0.897
Method:	Least Squares	F-statistic:	6.075e+04
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00
Time:	13:14:07	Log-Likelihood:	1.3460e+05
No. Observations:	35040	AIC:	-2.692e+05
Df Residuals:	35034	BIC:	-2.691e+05
Df Model:	5		
Covariance Type:	nonrobust		

#### Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
[2] The condition number is large, 1.48e+09. This might indicate that there are strong multicollinearity or other numerical problems.  
Final model R-squared: 0.8903953500042129  
Final model AIC: -209294.22094757787

#### Variance Inflation Factors for Final Model:

	Variable	VIF
0	const	186.898138
1	temperature_2m (°C)	1.285647
2	relative_humidity_2m (%)	1.377583
3	Lagging_Current_Reactive_Power_kVarh	1.377448
4	Lagging_Current_Power_Factor	1.049874
5	Leading_Current_Power_Factor	2.132558



\*\*More Model output and Interpretation on the next slide

## FINAL OLS Regression Model for CO<sub>2</sub>

### Without tempxhumid interaction and Leading Current Reactive Power

Final Model Summary:						
OLS Regression Results						
Dep. Variable:	CO2(tCO2)	R-squared:	0.897			
Model:	OLS	Adj. R-squared:	0.897			
Method:	Least Squares	F-statistic:	6.075e+04			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	1.3460e+05			
No. Observations:	35040	AIC:	-2.692e+05			
Df Residuals:	35034	BIC:	-2.691e+05			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.0397	0.800	-138.798	0.000	-0.040	-0.039
temperature_2m (°C)	6.051e-05	3.04e-06	19.880	0.000	5.45e-05	6.65e-05
relative_humidity_2m (%)	-2.463e-05	1.63e-06	-15.095	0.000	-2.78e-05	-2.14e-05
Lagging_Current_Reactive.Power_kVarh	0.0007	2.14e-06	321.334	0.000	0.001	0.001
Lagging_Current_Power_Factor	0.0004	2.05e-06	178.204	0.000	0.000	0.000
Leading_Current_Power_Factor	0.0002	1.32e-06	122.845	0.000	0.000	0.000
Omnibus:	4719.602	Durbin-Watson:	0.516			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	60217.016			
Skew:	0.137	Prob(JB):	0.00			
Kurtosis:	9.416	Cond. No.	1.43e+03			

#### Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.43e+03. This might indicate that there are strong multicollinearity or other numerical problems.

Final model R-squared: 0.8965953598042129

Final model AIC: -269194.21034717787

#### Variance Inflation Factors for Final Model:

	Variable	VIF
0	const	106.010198
1	temperature_2m (°C)	1.207647
2	relative_humidity_2m (%)	1.377503
3	Lagging_Current_Reactive.Power_kVarh	1.577468
4	Lagging_Current_Power_Factor	1.949874
5	Leading_Current_Power_Factor	2.111558

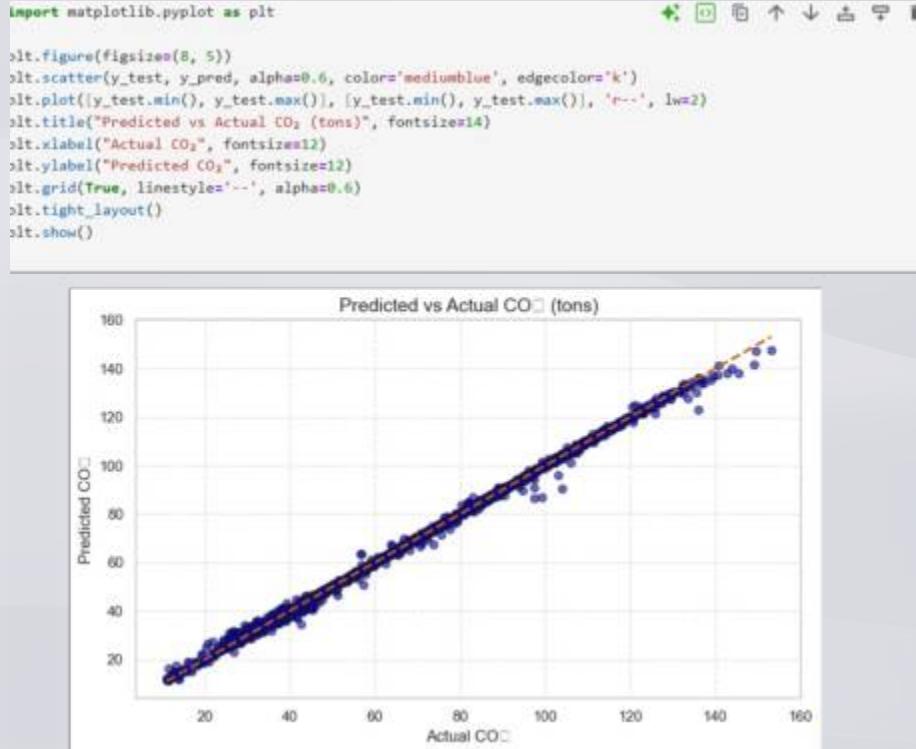
Model Stage	R <sup>2</sup>	AIC	Key Insight
Environment-only (No Interaction)	0.075	-192,440	Low explanatory power— temp & humidity contribute separately but weakly
Final CO <sub>2</sub> Model (w/o Interaction)	0.897	-269,194	Still excellent— power variables shoulder almost all explanatory weight
VIF Scores	< 2.2 for all terms	⚠️	Strong signal, low multicollinearity without the interaction term

## Interpretation

- Even without the temp × humidity interaction, environment-only model had minimal predictive strength.
- Adding power metrics, especially Lagging Reactive Power and both Power Factors, the model nearly recreates the full strength of the prior version with the interaction.
- Importantly, VIFs dropped, especially for temperature (from ~14.7 to ~1.2), meaning removing the interaction resolved multicollinearity cleanly without hurting model performance.
- The modeling suggest that for **CO<sub>2</sub> model**, humidity and temperature matter **primarily through their shared impact on operations or power efficiency**, already captured by the power variables.

# The OLS Model: Final Fit

## Simplicity Earned



While nonlinear models captured complexity across regimes and environmental layers, the final OLS regression delivered elegant simplicity. With well-behaved inputs and a refined feature set, the predicted vs. actual CO<sub>2</sub> plot shows near-perfect alignment with each point tightly clustered along the diagonal. This fit reflects not just statistical performance, but how far the modeling journey has come: from messy signals to interpretable clarity.

## OLS Regression Model Continues for Energy Usage kWh

### Without tempxhumid interaction

Environmental Variables Model (*without interaction of tempxhumid*):						
OLS Regression Results						
Dep. Variable:	Usage_kWh	R-squared:	0.078			
Model:	OLS	Adj. R-squared:	0.078			
Method:	Least Squares	F-statistic:	1473.			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	-1.7129e+05			
No. Observations:	35040	AIC:	3.426e+05			
Df Residuals:	35037	BIC:	3.426e+05			
Df Model:	2					
Covariance Type:	nonrobust					
<hr/>						
S	0.975]	coef	std err	t	P> t	[0.02
const	55.9434	0.632	88.493	0.000	54.70	
temperature_2m (°C)	0.4514	0.018	24.767	0.000	0.41	
relative_humidity_2m (%)	-0.4926	0.009	-53.852	0.000	-0.51	
1	-0.475					
Omnibus:	5405.050	Durbin-Watson:	0.191			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	8245.949			
Skew:	1.147	Prob(JB):	0.00			
Kurtosis:	3.620	Cond. No.	276.			
<hr/>						
Notes:						
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.						
R-squared: 0.07756887262874719						
AIC: 342587.64625590626						
Power variables found in dataframe: ['Lagging_Current_Reactive.Power_kVArh', 'Leading_Current_Reactive.Power_kVArh', 'Lagging_Current_Power_Factor', 'Leading_Current_Power_Factor']						
<hr/>						
Results of adding each power variable:						
Variable	AIC	R-squared	\			
0 Lagging_Current_Reactive.Power_kVArh	285387.734992	0.819715				
1 Leading_Current_Reactive.Power_kVArh	285364.892743	0.819843				
2 Lagging_Current_Power_Factor	265316.215698	0.898343				
3 Leading_Current_Power_Factor	260378.672452	0.911709				
<hr/>						
Coefficient	P-value					
0 1.797456	0.000000e+00					
1 0.057078	6.234469e-07					
2 0.691577	0.000000e+00					
3 0.397114	0.000000e+00					

Final Model Summary:						
OLS Regression Results						
Dep. Variable:	Usage_kWh	R-squared:	0.912			
Model:	OLS	Adj. R-squared:	0.912			
Method:	Least Squares	F-statistic:	6.829e+04			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	-1.3018e+05			
No. Observations:	35040	AIC:	2.604e+05			
Df Residuals:	35037	BIC:	2.604e+05			
Df Model:	6					
Covariance Type:	nonrobust					
<hr/>						
S	[0.825 0.975]	coef	std err	t	P> t	[0.02
const	-84.919	-82.110	-83.5150	0.717	-116.550	0.00
temperature_2m (°C)	0.090	0.113	0.1013	0.006	17.307	0.00
relative_humidity_2m (%)	-0.059	-0.047	-0.0533	0.003	-17.056	0.00
1	-0.323	0.411	1.459	1.475		
Lagging_Current_Reactive.Power_kVArh	0.323	0.411	0.3670	0.022	16.507	0.00
Leading_Current_Reactive.Power_kVArh	0.727	0.743	0.7349	0.004	185.417	0.00
Leading_Current_Power_Factor	0.386	0.408	0.3971	0.005	72.826	0.00

Final Model (Usage_kWh)						
(No Interaction + All Power Factors)						
• <b>R<sup>2</sup> = 0.912   AIC = 260,379</b>						

### Key Insights and Interpretation

Variable	R <sup>2</sup> ↑	AIC ↓	Coefficient	Interpretation
Lagging Current Reactive Power	0.820	285,388	1.80	Strong linear contribution to usage
Leading Reactive Power	0.820	285,365	0.06	Minimal but significant
Lagging Power Factor	0.898	265,316	0.69	Major driver of efficiency and demand
Leading Power Factor	0.912	260,379	0.40	Rounds out operational contribution

### Final Model (Usage\_kWh)

(No Interaction + All Power Factors)

• **R<sup>2</sup> = 0.912 | AIC = 260,379**

• Power factors dominate; temperature becomes secondary

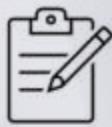
• All variables are **highly significant** (p < 0.001)

• **VIFs all < 10**, confirming multicollinearity is manageable

\*When environmental variables are modeled without interaction effects, they contribute modestly to energy demand. However, once operational power factors are included, explanatory strength soars (R<sup>2</sup> = 91.2%), affirming energy usage is driven more by internal system dynamics than by ambient conditions

## OLS Regression Model Continues for Energy Usage kWh

### Removing Leading\_Current\_Reactive\_Power\_kVarh variable



With the last model for Energy Usage kWh, I used the AIC for stepwise variable selection and surprisingly it kept the `Leading_Current_Reactive_Power_kVarh` variable. The Akaike Information Criterion (AIC) is a statistical measure that helps balance a model's goodness of fit with its complexity, guiding the selection of more parsimonious models. (Note, that the CO2 Model removed the variable. For this model, I am going to remove the variable above but still keep the AIC stepwise selection.)

Final Model Summary without 'Leading_Current_Reactive_Power_kVarh':					
OLS Regression Results					
Dep. Variable:	Usage_kWh	R-squared:	0.911		
Model:	OLS	Adj. R-squared:	0.911		
Method:	Least Squares	F-statistic:	7.174e+04		
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00		
Time:	13:14:07	Log-Likelihood:	-1.3032e+05		
No. Observations:	35040	AIC:	2.6066e+05		
Df Residuals:	35034	BIC:	2.607e+05		
Df Model:	5				
Covariance Type:	nonrobust				
<hr/>					
[0.025 0.975]		coef	std err	t	P> t
const	-76.943 -74.792	-75.8677	0.549	-138.252	0.00
temperature_2m (°C)	0.080 0.103	0.0914	0.006	15.635	0.00
relative_humidity_2m (%)	-0.060 -0.048	-0.0540	0.003	-17.213	0.00
Lagging_Current_Reactive.Power_kVarh	1.4585 1.467	1.4585	0.004	355.277	0.00
Lagging_Current_Power_Factor	0.7447 0.752	0.7447	0.004	189.335	0.00
Leading_Current_Power_Factor	0.3174 0.322	0.3174	0.003	124.818	0.00
<hr/>					
Omnibus:	5313.716	Durbin-Watson:	0.365		
Prob(Omnibus):	0.000	Jarque-Bera (JB):	17321.581		
Skew:	0.775	Prob(JB):	0.00		
Kurtosis:	6.076	Cond. No.	1.43e+03		

Notes:		
[1]	Standard Errors assume that the covariance matrix of the errors is correctly specified.	
[2]	The condition number is large, 1.43e+03. This might indicate that there are strong multicollinearity or other numerical problems.	
	Final model R-squared: 0.9110223395235306	
	Final model AIC: 260648.1427865061	
<hr/>		
Variance Inflation Factors for Final Model:		
	Variable	VIF
0	const	106.010198
1	temperature_2m (°C)	1.207647
2	relative_humidity_2m (%)	1.377503
3	Lagging_Current_Reactive.Power_kVarh	1.577468
4	Lagging_Current_Power_Factor	1.949874
5	Leading_Current_Power_Factor	2.111558

### Key Insights and Interpretation

Metric	With Variable	Without Variable
R-squared	0.9124	0.911
AIC	260,115	260,648
Coefficient (Removed Var)	+0.3667(significant)	✗ Removed
Multicollinearity Impact	VIF = 9.67(approaching threshold)	Lower overall VIFs
Interpretation	Slightly stronger fit, but added redundancy	Cleaner model with minimal loss of fit

Removing the `Leading_Current_Reactive_Power_kVarh` led to a minor dip in fit ( $\Delta R^2 = 0.0014$ ) and slight AIC increase, but significantly reduced complexity and potential collinearity. The decision emphasizes model parsimony while maintaining a high explanatory strength.



## NEXT STEPS for MODELS: Standardizing and Scaling Variable for Lasso and Cross -Validation

```
#print standardized predictors
print(X_scaled_df.head())
# Compare statistics
print("Unscaled Variables")
df[['temperature_2m (°C)', 'relative_humidity_2m (%)',
   'Lagging_Current_Reactive.Power_kVarh', 'Lagging_Current_Power_Factor',
   'Leading_Current_Power_Factor']].head()
```

	temperature_2m (°C)	relative_humidity_2m (%)	Lagging_Current_Reactive.Power_kVarh	Lagging_Current_Power_Factor	Leading_Current_Power_Factor
0	-1.512800	-0.388125	-0.428916	-0.595418	0.513264
1	-1.512800	-0.388125	-0.525931	-0.729772	0.513268
2	-1.512800	-0.388125	-0.596279	-0.548264	0.513268
3	-1.512816	-0.288913	-0.581106	-0.608009	0.513268
4	-1.512816	-0.288913	-0.523458	-0.858117	0.513268

Unscaled Variables

	temperature_2m (°C)	relative_humidity_2m (%)	Lagging_Current_Reactive.Power_kVarh	Lagging_Current_Power_Factor	Leading_Current_Power_Factor
0	-1.3	63.0	2.85	73.21	100.0
1	-1.3	63.0	4.40	66.77	100.0
2	-1.3	63.0	3.28	70.28	100.0
3	-1.5	65.0	3.56	68.00	100.0
4	-1.5	65.0	4.50	64.72	100.0

Good results on models but still shows potential multicollinearity. Next step is to run Lasso and Cross-Validation with other potential metrics. However, need to standardize variables

Standardization /scaling is a preprocessing step that prepares data for modeling.

- **LASSO penalizes based on coefficient size**, so a variable with larger units (e.g., kWh) could be unfairly penalized unless scaled
- **Cross-validation performance metrics** (like RMSE or MAE) will be more interpretable across folds
- **Multicollinearity effects** (e.g., VIF) also become easier to identify when all predictors are mean-centered and scaled. This ensures all subsequent analyses, (cross-validation, LASSO, metrics like BIC), use the same scaled data, avoiding inconsistencies.

# CROSS -Validation

```
#Checking and picking my Models (stable?)  
# Cross-validation for CO2 model  
model = LinearRegression()  
cv_scores_co2 = cross_val_score(model, X_scaled_df, y_co2, cv=5, scoring='neg_mse')  
mse_co2 = -cv_scores_co2.mean()  
rmse_co2 = np.sqrt(mse_co2)  
print(f"CO2 Model - 5-Fold CV MSE: {mse_co2:.6f}, RMSE: {rmse_co2:.6f}")  
  
# Cross-validation for kWh model  
cv_scores_kwh = cross_val_score(model, X_scaled_df, y_kwh, cv=5, scoring='neg_mse')  
mse_kwh = -cv_scores_kwh.mean()  
rmse_kwh = np.sqrt(mse_kwh)  
print(f"kWh Model - 5-Fold CV MSE: {mse_kwh:.6f}, RMSE: {rmse_kwh:.6f}")  
  
# Test model without Leading_Current_Power_Factor.  
X_reduced = X_scaled_df.drop(columns=['Leading_Current_Power_Factor'])  
cv_scores_reduced = cross_val_score(model, X_reduced, y_kwh, cv=5, scoring='neg_mse')  
# Fixed: Calculate mean first, then negate  
mse_reduced = -cv_scores_reduced.mean() # Just take the mean of negative scores  
rmse_reduced = np.sqrt(mse_reduced)  
print(f"kWh Model (without Leading_Current_Power_Factor) - CV MSE: {mse_reduced:.6f}, RMSE: {rmse_reduced:.6f}")  
  
CO2 Model - 5-Fold CV MSE: 0.000031, RMSE: 0.005532  
kWh Model - 5-Fold CV MSE: 115.588644, RMSE: 10.751216  
kWh Model (without Leading_Current_Power_Factor) - CV MSE: 155.424307, RMSE: 12.466929
```



## Cross-Validation Results Summary

Model	CV MSE	CV RMSE	Interpretation
CO <sub>2</sub> Model	0.000031	0.0055	Incredibly tight error and model generalizes well on CO <sub>2</sub> prediction
kWh Model (full)	115.59	10.75	Solid performance and low error given operational variability
kWh Model (minus Leading PF)	155.42	12.47	Error rises sharply without Leading PF and confirms its explanatory power

# Plot of Lasso's Coefficients for both CO<sub>2</sub> and Energy Usage kWh

**Overview:** Where energy demand reflects an intricate balance of power quality and environmental conditions, CO<sub>2</sub> emissions reduce to a single operational driver, reactive power. The LASSO coefficient comparison show their structural differences.

```
#Coefficient Comparison Charts
#Using a side-by-side bar plot to show scaled LASSO coefficients for both CO2 and kWh models
import matplotlib.pyplot as plt
import pandas as pd

# Assuming lasso_co2 and lasso_kwh are already defined and fitted

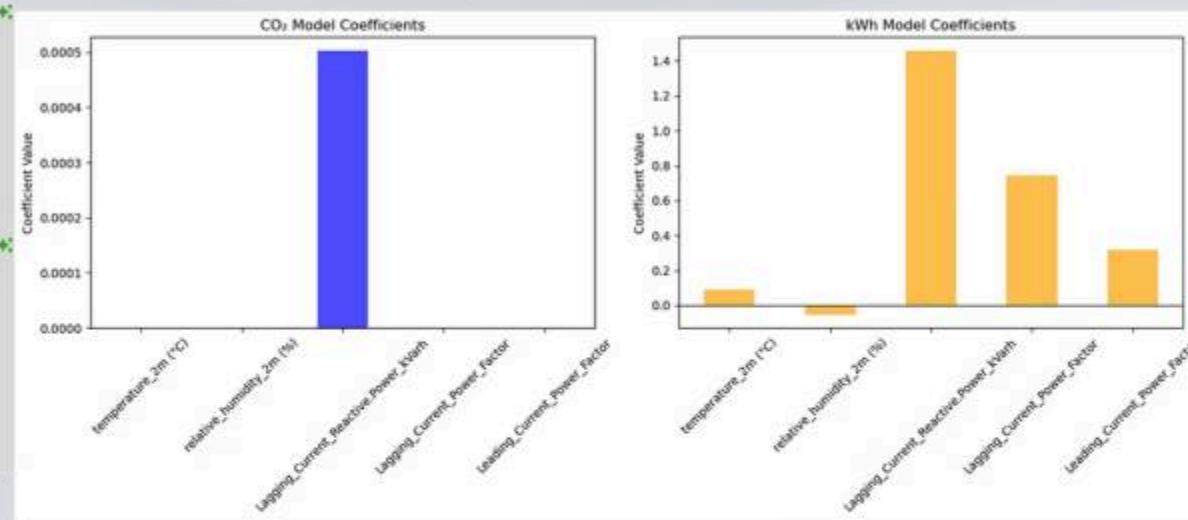
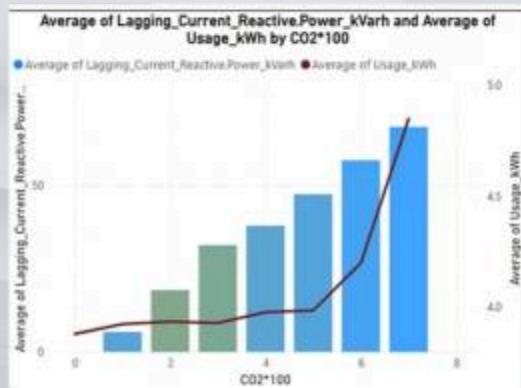
# Create the coefficient comparison dataframe
coefs_df = pd.DataFrame({
    'y_co2': lasso_co2.coef_,
    'y_kwh': lasso_kwh.coef_
}, index=predictors)

#Coefficient Comparison Charts
#Using a side-by-side bar plot to show scaled LASSO coefficients for both CO2 and kWh models
import matplotlib.pyplot as plt
import pandas as pd

# Assuming lasso_co2 and lasso_kwh are already defined and fitted

# Create the coefficient comparison dataframe
coefs_df = pd.DataFrame({
    'y_co2': lasso_co2.coef_,
    'y_kwh': lasso_kwh.coef_
}, index=predictors)
```

Previous Comparison of Lagging Reactive Power, CO<sub>2</sub> and Energy Usage kWh



## CO<sub>2</sub> Model

- Only Lagging Reactive Power survives LASSO penalty.
- That single non-zero bar is emissions proxy: if it draws VARh, it emits CO<sub>2</sub>.
- All other coefficients are shrunk to zero, including both power factors and environmental terms.

## kWh Model (Energy Usage)

- Lagging Reactive Power towers above with the strongest signal.
- Lagging Power Factor follows, meaning efficiency changes directly affect demand.
- Leading Power Factor and environmental variables (temp, humidity) contribute modestly.
- This confirms that the Usage\_kWh is operationally driven, but responsive to environment at the margins

# Ordinary Least Squares (OLS) Model with Lasso Selected Variables

```
# Fit OLS with LASSO-selected variables (Ordinary Least Squares)
#Filtering out features where Lasso coefficients are zero, focus on important predictors
lasso_selected_co2 = X_scaled_sm.loc[:, ['const'] + [col for col, coef in zip(predictors, lasso_co2.coef_) if coef != 0]]
lasso_selected_kwh = X_scaled_sm.loc[:, ['const'] + [col for col, coef in zip(predictors, lasso_kwh.coef_) if coef != 0]]

model_co2_lasso = sm.OLS(y_co2, lasso_selected_co2).fit()
model_kwh_lasso = sm.OLS(y_kwh, lasso_selected_kwh).fit()

print("\nCO2 Model (LASSO-selected variables):")
print(model_co2_lasso.summary())
print("\nkWh Model (LASSO-selected variables):")
print(model_kwh_lasso.summary())
```

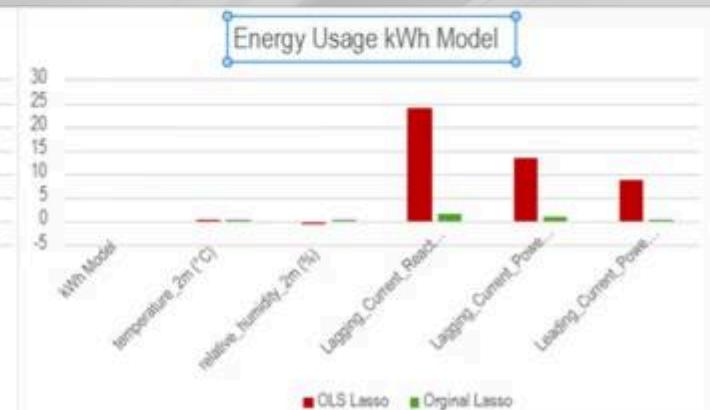
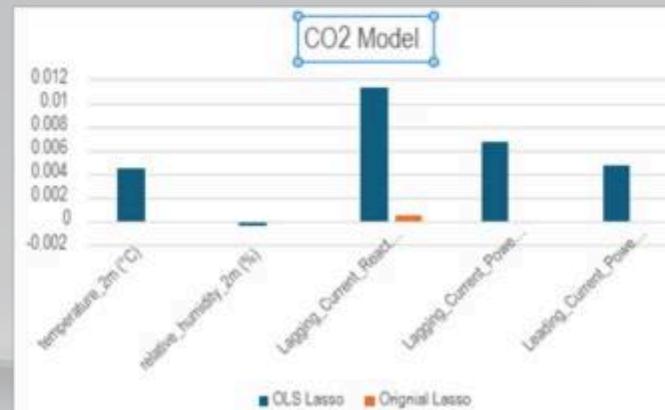
## Coefficients and Comparison from Original Lasso

### CO2 Model - LASSO Coefficients:

```
temperature_2m (°C)          0.000444
relative_humidity_2m (%)      -0.000378
Lagging_Current_Reactive.Power_kVArh 0.011258
Lagging_Current_Power_Factor    0.006716
Leading_Current_Power_Factor   0.004708
dtype: float64
Optimal alpha: 0.000094
```

### kWh Model - LASSO Coefficients:

```
temperature_2m (°C)          0.280719
relative_humidity_2m (%)      -0.630660
Lagging_Current_Reactive.Power_kVArh 24.016861
Lagging_Current_Power_Factor   13.353312
Leading_Current_Power_Factor   8.706943
dtype: float64
Optimal alpha: 0.369493
```



## Ordinary Least Squares (OLS) Model with Lasso Selected Variables

### Model outputs and Evaluation

LASSO models maintain high R<sup>2</sup> (CO<sub>2</sub> = 0.897, kWh = 0.911) Cross-validation RMSEs remain tight (5-fold cross-validation)

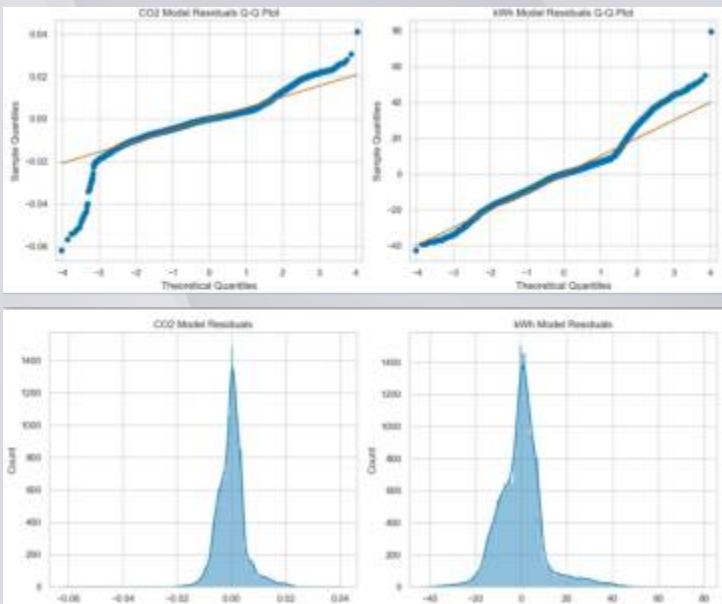
CO2 Model (LASSO-selected variables): OLS Regression Results							kWh Model (LASSO-selected variables): OLS Regression Results						
Dep. Variable:	CO2(tCO2)	R-squared:	0.897	Dep. Variable:	Usage_kWh	R-squared:	0.911						
Model:	OLS	Adj. R-squared:	0.897	Model:	OLS	Adj. R-squared:	0.911						
Method:	Least Squares	F-statistic:	6.075e+04	Method:	Least Squares	F-statistic:	7.174e+04						
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00	Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00						
Time:	13:14:10	Log-Likelihood:	1.3460e+05	Time:	13:14:10	Log-Likelihood:	-1.3032e+05						
No. Observations:	35040	AIC:	-2.692e+05	No. Observations:	35040	AIC:	2.606e+05						
Df Residuals:	35034	BIC:	-2.691e+05	Df Residuals:	35034	BIC:	2.607e+05						
Df Model:	5			Df Model:	5								
Covariance Type:	nonrobust			Covariance Type:	nonrobust								
	coef	std err.	t	P> t	[0.025	0.975]		coef	std err.	t	P> t	[0.025	0.975]
const	0.0115	2.77e-05	415.336	0.000	0.011	0.012	const	27.3869	0.053	\$13.883	0.000	27.282	27.491
temperature_2m (°C)	0.0086	3.85e-05	19.888	0.000	0.001	0.001	temperature_2m (°C)	0.9157	0.059	15.635	0.000	0.801	1.031
relative_humidity_2m (%)	-0.0005	3.26e-05	-15.095	0.000	-0.001	-0.000	relative_humidity_2m (%)	-1.0768	0.063	-17.213	0.000	-1.199	-0.954
lagging_Current_Reactive.Power_kVarh	0.0112	3.48e-05	321.334	0.000	0.011	0.011	lagging_Current_Reactive.Power_kVarh	23.7826	0.067	355.277	0.000	23.651	23.914
lagging_Current_Power_Factor	0.0069	3.87e-05	170.204	0.000	0.007	0.007	lagging_Current_Power_Factor	14.0911	0.074	189.335	0.000	13.945	14.237
leading_Current_Power_Factor	0.0050	4.03e-05	122.845	0.000	0.005	0.005	leading_Current_Power_Factor	9.6670	0.077	124.818	0.000	9.515	9.819
Omnibus:	4719.602	Durbin-Watson:	0.516		Omnibus:	5313.716	Durbin-Watson:	0.365					
Prob(Omnibus):	0.000	Jarque-Bera (JB):	60217.016		Prob(Omnibus):	0.000	Jarque-Bera (JB):	17321.581					
Skew:	0.137	Prob(JB):	0.00		Skew:	0.775	Prob(JB):	0.00					
Kurtosis:	9.416	Cond. No.	2.64		Kurtosis:	6.076	Cond. No.	2.64					

- It appears that CO<sub>2</sub> emissions are operationally lean?
- LASSO compresses some coefficients, but reactive energy, lagging power factor, and leading power factor stand out and yields energy-derived emissions structure.

- kWh usage shows high coefficients across operational and environmental variables.
- This appears to reflect broader sensitivity to real-world load and ambient conditions.

\*High Jarque-Bera (JB) Test Concerns (Normality of Residual) possibly due to Outliers, Heavy-tailed errors, omitted interactions and Autocorrelation (Time-Based Dependence). Will try to address in other models.

# RESIDUAL PLOTS FOR LASSO OLS MODEL



## CO<sub>2</sub> Model Residuals (Left Panels)

### Top-Left: Q-Q Plot

- Exhibits **subtle curvature**, especially in the **lower region**, where residuals drop markedly beneath the fitted line.
- The upper region remains relatively aligned, but begins to **veer away** toward higher predicted CO<sub>2</sub>, suggesting **nonlinearity concentrated in low-to-mid ranges**.

## Energy Usage (kWh) Residuals (Right Panels)

### Top-Right: Q-Q Plot

- Closely tracks the **fitted curve** until the mid-upper range, where residuals begin to arc upward which is a possible sign of **saturation effects** at higher energy levels or not specified load interactions.

## CO<sub>2</sub> Model Residuals (Left Panels)

- The **extended left tail** likely corresponds to **zero-inflated emissions**, e.g., operational states with minimal output or downtime and **zero cluster (transform and imputation needed)**
- This tail skewness is also **statistically reflected in**:
  - Skew ~ +0.14**, slightly right-leaning overall but influenced by extreme low-end residuals
  - Kurtosis ~ 9.4**, indicating tight central clustering with heavy tails
  - Jarque-Bera statistic > 60,000**, rejecting normality and flagging non-symmetric error behavior

## Energy Usage (kWh) Residuals (Left Panels)

- Smaller left tail and a larger, more dramatic right tail** is suggesting more extreme over-predictions of energy usage at peak levels.
- Combined with **moderate kurtosis (~6.1)** and **strong JB test**, this hints at skew and heavy tails but slightly less severity than in the CO<sub>2</sub> model.

# Continue Testing for Non-Normality and Skewness

Shapiro-Wilk Test for CO2(tCO2): Statistic=0.7307, p-value=0.0000

Shapiro-Wilk Test for Usage\_kWh: Statistic=0.7532, p-value=0.0000

Shapiro-Wilk Test for Lagging\_Current\_Reactive.Power\_kVarh: Statistic=0.7686, p-value=0.0000

CO2 Skewness: 1.1494

Usage\_kWh Skewness: 1.1974

## Shapiro-Wilk Results

Variable	W Statistic	p-value	Conclusion
CO <sub>2</sub> (tCO <sub>2</sub> )	0.7307	0.0000	✗ Not normally distributed
Usage_kWh	0.7532	0.0000	✗ Not normal
Lagging Reactive Power (kVArh)	0.7686	0.0000	✗ Not normal

## Skewness Check

Variable	Skewness	Interpretation
CO <sub>2</sub> (tCO <sub>2</sub> )	+1.15	Right-skewed → long tail of higher emissions
Usage_kWh	+1.20	Right-skewed → some extreme usage spikes

Both CO<sub>2</sub> and energy usage are significantly right-skewed and non-normal, as confirmed by Shapiro-Wilk tests and skewness metrics. These distributional features justify our use of log transformation and penalized regression for improved fit and interpretability.

# Exploring Data Transformations: Strengthening Model Validity

To further improve model reliability and meet regression assumptions, I explored a series of **transformations** aimed at reducing skewness, stabilizing variance, and enhancing residual symmetry:

- **Log Transformation (Previous Models):** Ideal for positively skewed variables like energy usage and CO<sub>2</sub> emissions. It helped reduce error inflation and improve normality, especially useful for interpretability when modeling proportional or multiplicative relationships.
- **Yeo-Johnson Transformation:** Chosen for its flexibility. Unlike log, it handles **zero and negative values** reliably. This was especially valuable for the CO<sub>2</sub> data, where zeros and low-load states were common. Yeo-Johnson substantially improved model fit and residual distribution.
- **Quantile Transformation (Normal Output):** Rescales the data to follow a **standard normal distribution** using rank mapping. While it offers symmetry, it may reduce interpretability and distort scale-sensitive effects.

These transformed models offer **alternate perspectives on the same relationships**, helping validate possible conclusions while pushing closer to statistical rigor. Each method balances **fit, normality, and interpretability** differently, while offering a diverse toolkit for predictive modeling and storytelling.

# Modeling Enhancements: Exploring Log Transformation



To improve residual behavior and increase confidence in model assumptions, I applied a log transformation to both CO<sub>2</sub> emissions and energy usage (kWh). The upcoming slide revisits the Usage\_kWh model, which was the first model I built, now with transformed values. While I was initially hesitant about applying the same transformation to CO<sub>2</sub>, due to its zero-inflation suspect, I couldn't resist testing what impact it might have on residual symmetry, skewness, and information criteria.

## Log Transform OLS Regression Model - CO2

CO2 Model (Log-transformed): OLS Regression Results						
Dep. Variable:	log_CO2	R-squared:	0.898			
Model:	OLS	Adj. R-squared:	0.898			
Method:	Least Squares	F-statistic:	6.188e+04			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:12	Log-Likelihood:	1.3561e+05			
No. Observations:	350408	AIC:	-2.712e+05			
Df Residuals:	35034	BIC:	-2.712e+05			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	0.0113	2.7e-05	420.289	0.000	0.011	0.011
temperature_2m (°C)	0.0006	2.96e-05	20.851	0.000	0.001	0.001
relative_humidity_2m (%)	-0.0005	3.16e-05	-15.256	0.000	-0.001	-0.000
Lagging_Current_Reactive.Power_kVarh	0.0109	3.39e-05	322.642	0.000	0.011	0.011
Lagging_Current_Power_Factor	0.0061	3.77e-05	181.473	0.000	0.007	0.007
Leading_Current_Power_Factor	0.0049	3.92e-05	125.792	0.000	0.005	0.005
Omnibus:	4692.239	Durbin-Watson:	0.511			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	60845.917			
Skew:	0.104	Prob(JB):	0.00			
Kurtosis:	9.452	Cond. No.	2.64			

### Result and Interpretation of Model

Metric	Original CO <sub>2</sub>	Log CO <sub>2</sub>	Δ (Improvement)
Adjusted R <sup>2</sup>	0.8966	0.8983	+0.0017 (negligible)
BIC	-368,582.64	<b>-370,589.95</b>	+ 2,007 points <span style="color: green;">(+) 2,007 points</span>
Skewness	~+1.15	<b>0.1</b>	<span style="color: green;">✓</span> Vastly reduced
Kurtosis	~9.42	~9.45	No meaningful change
Residual Plots	Moderate non-normality	Slightly improved symmetry	Mildly better

Although the log transformation did not significantly improve R<sup>2</sup>, it reduced skewness and improved distributional symmetry, leading to a 2,000-point drop in BIC. While optional, it strengthens the statistical assumptions behind inference and may benefit future forecasting or time series modeling.

## Log Transform OLS Regression Model - Energy Usage kWh

```
# Re-run Ordinary Least Squares(OLS) with log-transformed kWh
model_kwh_log = sm.OLS(df['log_kwh'], lasso_selected_kwh).fit()
print("\n{kwh Model (Log-transformed):")
print(model_kwh_log.summary())

# Calculate BIC for log-transformed model.
def calculate_bic(model, n, y, X):
    mse = mean_squared_error(y, model.fittedvalues)
    k = len(model.params)
    bic = n * np.log(mse) + k * np.log(n)
    return bic

bic_kwh_log = calculate_bic(model_kwh_log, len(df['log_kwh']), df['log_kwh'], lasso_selected_kwh)
print(f"\nLog kWh Model BIC: {bic_kwh_log:.2f}, Adjusted R-squared: {model_kwh_log.rsquared_adj:.4f}")
print("Original kWh BIC: 161259.72, Adjusted R-squared: 0.9118")
```

### Result and Interpretation of Model

Model	Adjusted R <sup>2</sup>	BIC	Interpretation
Original kWh	0.911	161,259.72	Solid fit, but residuals showed skew/kurtosis
Log-transformed kWh	0.9345	-77,984.57	💡 Substantially better fit + reduced error structure

Since BIC penalizes for number of predictors; I assume that this drop isn't just due to added complexity. I believe it's a clear signal that the log transformation corrected skew and stabilized variance, aligning better with linear model assumptions.

kWh Model (Log-transformed):

OLS Regression Results

Dep. Variable:	log_kwh	R-squared:	0.934			
Model:	OLS	Adj. R-squared:	0.934			
Method:	Least Squares	F-statistic:	9.995e+04			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:12	Log-Likelihood:	-10696.			
No. Observations:	35040	AIC:	2.140e+04			
Df Residuals:	35034	BIC:	2.145e+04			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	2.5467	0.002	1451.768	0.000	2.543	2.550
temperature_2m (°C)	0.0836	0.002	43.350	0.000	0.080	0.087
relative_humidity_2m (%)	-0.0389	0.002	-18.913	0.000	-0.043	-0.035
Lagging_Current_Reactive.Power_kVarh	0.6641	0.002	301.430	0.000	0.660	0.668
Lagging_Current_Power_Factor	0.8350	0.002	340.878	0.000	0.830	0.840
Leading_Current_Power_Factor	0.6507	0.003	255.258	0.000	0.646	0.656
Omnibus:	3568.753	Durbin-Watson:	0.308			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	5458.295			
Skew:	-0.763	Prob(JB):	0.00			
Kurtosis:	4.187	Cond. No.	2.64			

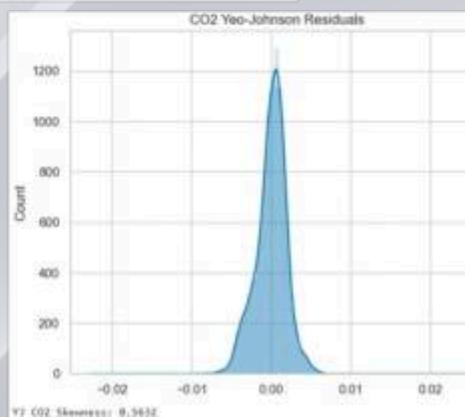
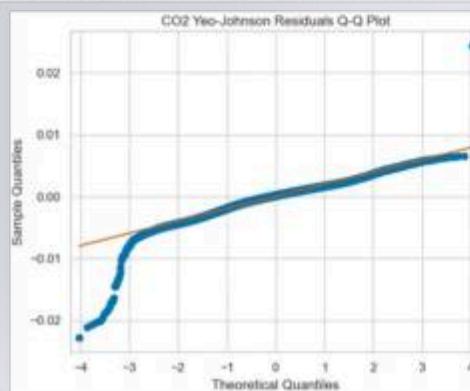
Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Log kWh Model BIC: -77984.57, Adjusted R-squared: 0.9345  
Original kWh BIC: 161259.72, Adjusted R-squared: 0.9118

# Yeo-Johnson Model with CO<sub>2</sub>

```
CO2 Model (Yeo-Johnson):
OLS Regression Results
*****
Dep. Variable: yj_CO2 R-squared: 0.914
Model: OLS Adj. R-squared: 0.914
Method: Least Squares F-statistic: 7.403e+04
Date: Mon, 30 Jun 2025 Prob (F-statistic): 0.00
Time: 13:14:13 Log-Likelihood: 1.6838e+05
No. Observations: 35040 AIC: -3.368e+05
Df Residuals: 35034 BIC: -3.367e+05
Df Model: 5
Covariance Type: nonrobust
*****
            coef    std err      t      P>|t|      [0.025    0.975]
-----
const        0.0053  1.06e-05  504.960   0.000     0.005     0.005
temperature_2m (°C)  0.0007  1.16e-05   56.888   0.000     0.001     0.001
relative_humidity_2m (%) -0.0002  1.24e-05  -17.066   0.000    -0.000    -0.000
Lagging_Current_Reactive_Power_kVarh  0.0036  1.33e-05  270.855   0.000     0.004     0.004
Lagging_Current_Power_Factor  0.0040  1.48e-05  272.082   0.000     0.004     0.004
Leading_Current_Power_Factor  0.0034  1.54e-05  219.589   0.000     0.003     0.003
-----
Omnibus: 7203.426 Durbin-Watson: 0.450
Prob(Omnibus): 0.000 Jarque-Bera (JB): 48850.466
Skew: -0.820 Prob(JB): 0.00
Kurtosis: 8.547 Cond. No. 2.64
*****
```



## Result and Interpretation of Model

Metric	Value	Interpretation
Adjusted R <sup>2</sup>	0.914	✓ Excellent fit (strongest yet)
BIC	-336,700	Significant drop from log model
Skewness	-0.82	Left-leaning tail (zero emission states)
Kurtosis	8.55	Tighter core with heavy tails
JB Stat	~48,850	⚠ Non-normal residuals persist

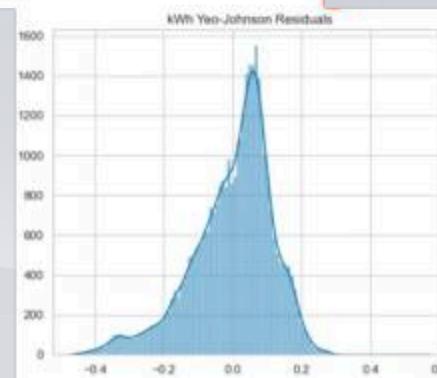
- The Yeo-Johnson transformation **dampens extreme skewness** while preserving linear relationships, especially around zero-heavy data.
- Residual symmetry improves though **left skew remains**, likely tied to **zero or low-emission operating states**.
- Compared to the log-transformed model, you gain:
  - Better R<sup>2</sup>**
  - Lower BIC**
  - More centered alignment**
- Concern** -The Q-Q plot looks great, except the large curvature at start of plot

# Yeo-Johnson Model with Energy Usage kWh

kWh Model (Yeo-Johnson):

OLS Regression Results

Dep. Variable:	yj_kwh	R-squared:	0.915			
Model:	OLS	Adj. R-squared:	0.915			
Method:	Least Squares	F-statistic:	7.533e+04			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	13:14:14	Log-Likelihood:	-25342.			
No. Observations:	35040	AIC:	-5.067e+04			
Df Residuals:	35034	BIC:	-5.062e+04			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	1.4369	0.001	2290.957	0.000	1.436	1.438
temperature_2m (°C)	0.0233	0.001	33.780	0.000	0.022	0.025
relative_humidity_2m (%)	-0.0111	0.001	-15.085	0.000	-0.013	-0.010
Lagging_Current_Reactive.Power_kVArh	0.1791	0.001	227.348	0.000	0.178	0.181
Lagging_Current_Power_Factor	0.2889	0.001	329.848	0.000	0.287	0.291
Leading_Current_Power_Factor	0.2229	0.001	244.597	0.000	0.221	0.225
Omnibus:	3550.102	Durbin-Watson:	0.333			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	4869.116			
Skew:	-0.823	Prob(JB):	0.00			
Kurtosis:	3.789	Cond. No.	2.64			



## Result and Interpretation of Model

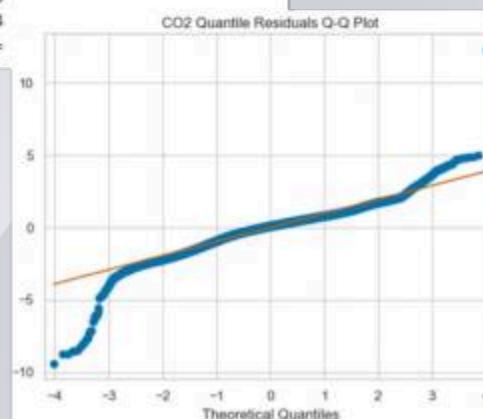
Metric	Value	Interpretation
Adjusted R <sup>2</sup>	0.915	✓ High explanatory power
BIC	-50,620	Substantial improvement from untransformed model
Skewness	-0.82	Left-skewed residuals with likely over suppression in low-usage cases
Kurtosis	3.79	Closer to normal with improved tail
Jarque-Bera	4,869.12	✗ Residuals still non-normal but improved

- The Yeo-Johnson transformation **improved model quality** by addressing skewness and stabilizing variance, especially in low-energy usage zones.
- Unlike log, this method handled **zero and near-zero values** gracefully without compressing interpretability.
- Residuals are **more symmetrical and less heteroskedastic**, but residuals plot still shows curvature.

# Quantile Model with CO<sub>2</sub>

```
CO2 Model (Quantile):
OLS Regression Results
=====
Dep. Variable: qt_CO2 R-squared: 0.897
Model: OLS Adj. R-squared: 0.897
Method: Least Squares F-statistic: 6.124e+04
Date: Mon, 30 Jun 2025 Prob (F-statistic): 0.00
Time: 13:14:15 Log-Likelihood: -48719.
No. Observations: 35048 AIC: 9.745e+04
Df Residuals: 35034 BIC: 9.750e+04
Df Model: 5
Covariance Type: nonrobust
=====
            coef    std err      t   P>|t|    [ 0.825    0.975 ]
const        -2.7347   0.005  -526.696   0.000    -2.745    -2.725
temperature_2m (°C)  0.3525   0.006   61.778   0.000    0.341     0.364
relative_humidity_2m (%) -0.0922   0.006  -15.137   0.000    -0.104    -0.080
Lagging_Current_Reactive.Power_kVarh  1.3537   0.007  207.582   0.000    1.341     1.366
Lagging_Current_Power_Factor  2.0292   0.007  279.878   0.000    2.015     2.043
Leading_Current_Power_Factor  1.7591   0.008  233.146   0.000    1.744     1.774
=====
Omnibus: 5410.027 Durbin-Watson: 0.477
Prob(Omnibus): 0.000 Jarque-Bera (JB): 29138.357
Skew: -0.640 Prob(JB): 0.00
Kurtosis: 7.280 Cond. No. 2.64
=====
```

Quantile transformation maps CO<sub>2</sub> outputs to a standard normal distribution, improving symmetry and residual alignment. Though less interpretable, the model maintains strong predictive power and helps verify the robustness of our variable relationships.



## Result and Interpretation of Model

Metric	Value	Interpretation
Adjusted R <sup>2</sup>	0.897	✓ Predictive strength sustained
BIC	~97,500	Higher than log and Yeo-Johnson models
Skewness	-0.64	Left tail elongation; tighter than raw, not perfect
Kurtosis	7.28	Still peaked with heavy central concentration
JB Test (Stat)	~29,138	✗ Non-normal residuals persist

- The quantile transformation forces the output into a normal distribution by ranking and stretching values along a Gaussian curve (great for fixing skew) but may distort relationships.
- Predictive strength is preserved, but the coefficients are harder to interpret.
- Residuals still show curvature at top and bottom.
- The tail behavior (Jarque-Bera, kurtosis) still flags potential outlier effects.

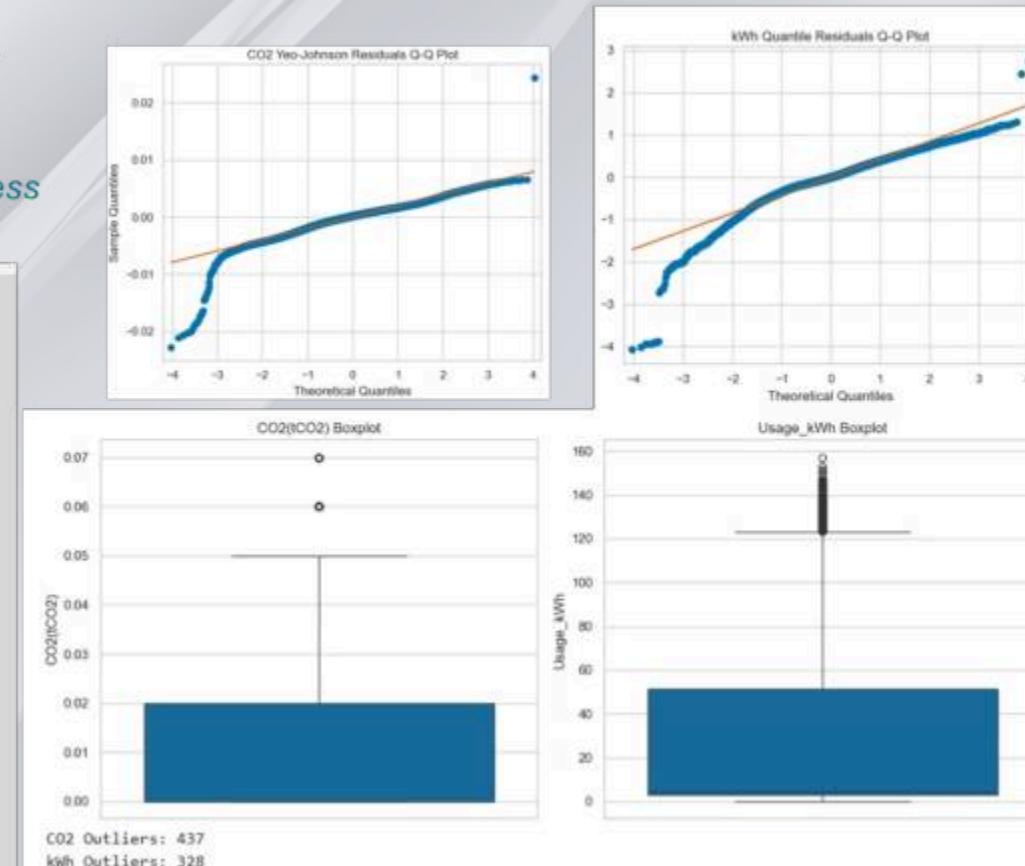
## Checking for outliers for both CO<sub>2</sub> and Energy Usage kWh

After a closer review of residuals from my previous models, I noticed several patterns that raise concerns about the overall fit, or lack thereof, in certain areas. I suspect that some of these discrepancies may stem from underlying outliers, and will investigate and address them as part of the model refinement process

```
# Checking for outliers which could be the problem with normal fit
#Boxplots for outliers
plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
sns.boxplot(y=df['CO2(tCO2)'])
plt.title('CO2(tCO2) Boxplot')
plt.subplot(1, 2, 2)
sns.boxplot(y=df['Usage_kWh'])
plt.title('Usage_kWh Boxplot')
plt.tight_layout()
plt.show()

# Identify outliers (IQR method)
def detect_outliers(series):
    Q1, Q3 = series.quantile([0.25, 0.75])
    IQR = Q3 - Q1
    lower, upper = Q1 - 1.5 * IQR, Q3 + 1.5 * IQR
    return series[(series < lower) | (series > upper)]

co2_outliers = detect_outliers(df['CO2(tCO2)'])
kwh_outliers = detect_outliers(df['Usage_kWh'])
print(f"CO2 Outliers: {len(co2_outliers)}")
print(f"kWh Outliers: {len(kwh_outliers)}")
```



# Model Comparison: Outlier Removal Effects

## Revisiting the Yeo-Johnson with No Outliers

```
# Assuming df, lasso_selected_co2, lasso_selected_kwh from previous
#Remove outliers
def remove_outliers(series):
    Q1, Q3 = series.quantile([0.25, 0.75])
    IQR = Q3 - Q1
    lower, upper = Q1 - 1.5 * IQR, Q3 + 1.5 * IQR
    return (series >= lower) & (series <= upper)

# Filter outliers
mask_co2 = remove_outliers(df['CO2(tCO2)'])

mask_kwh = remove_outliers(df['Usage_kWh'])
df_no_outliers = df[mask_co2 & mask_kwh].copy()

#Re-apply Yeo-Johnson which had the strongest model so far. will add info later
from sklearn.preprocessing import PowerTransformer
pt = PowerTransformer(method='yeo-johnson', standardize=False)
df_no_outliers['yj_CO2'] = pt.fit_transform(df_no_outliers[['CO2(tCO2)']])
df_no_outliers['yj_kWh'] = pt.fit_transform(df_no_outliers[['Usage_kWh']])
```

Metric	CO <sub>2</sub> No Outliers	CO <sub>2</sub> All Data	kWh No Outliers	kWh All Data
Adjusted R <sup>2</sup>	0.9142	0.9135	0.9139	0.9149
BIC	-434668.76	-436144.55	-153391.67	-150060.26
Skewness	0.5771	0.5632	0.3629	0.3442
Kurtosis	1.4722	—	1.3013	—

## Key Points

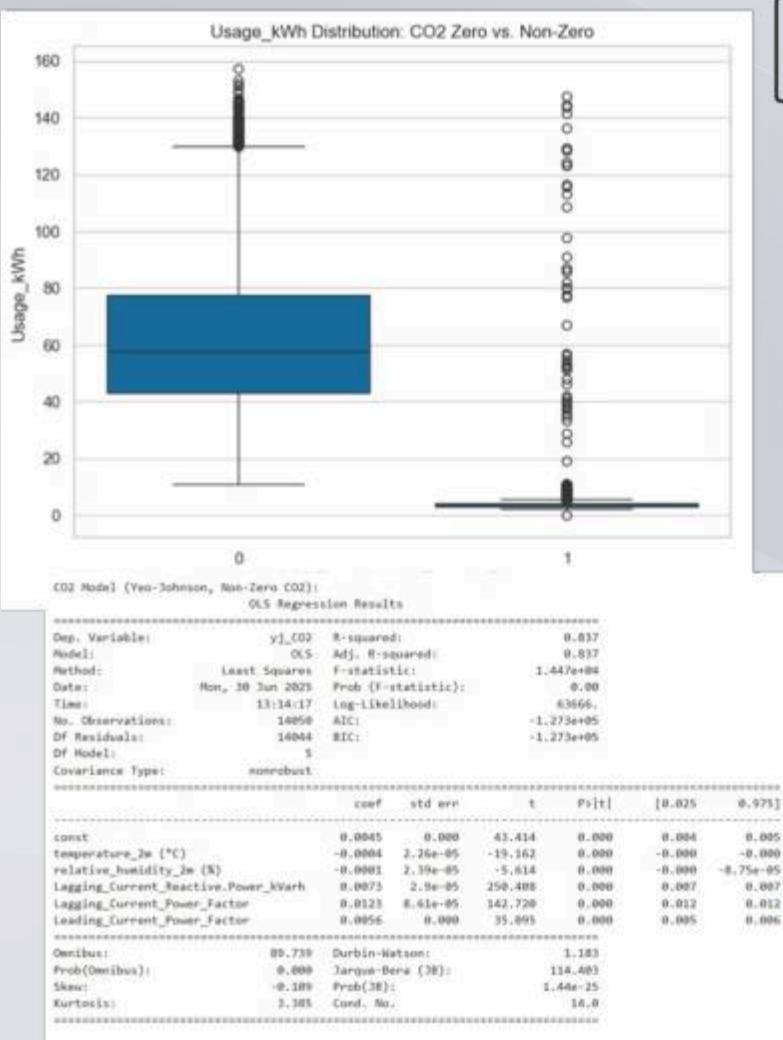
- BIC favors the full data, nudging toward keeping all points (even the quirky ones).
- Adjusted R-squared barely shifts, suggesting the models are resilient either way.
- Skewness and kurtosis steadily improve after trimming , especially in CO<sub>2</sub>, which went from moderately skewed to more symmetric.

### TAKEAWAY:

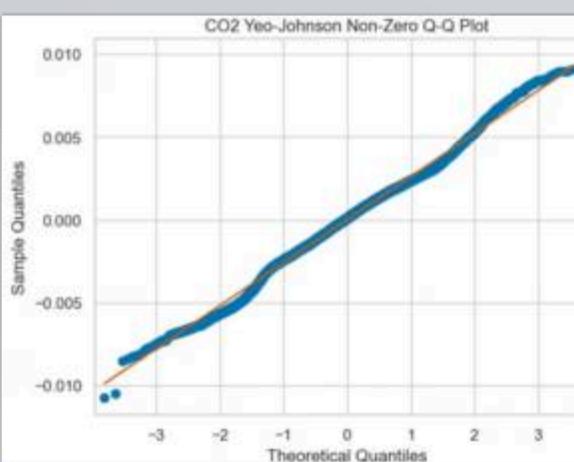
Outlier removal reduced tail risk (skewness, kurtosis) with minimal R<sup>2</sup> loss, but higher BIC suggests a trade-off in model parsimony.

## CO<sub>2</sub> Non-Zero Model Analysis

BOXPLOT: CO<sub>2</sub> Non-Zero and Zero Comparison



To better understand the structural behavior of CO<sub>2</sub> emissions, I ran an OLS model using only non-zero emission records. While this filtering reduces the overall data volume and excludes some genuinely meaningful zeros, the goal was to examine how the continuous emission patterns behave in isolation. By excluding zeros, I anticipated a modest decline in model fit metrics (e.g., adjusted R<sup>2</sup>) but hoped to improve distributional traits like skewness, kurtosis, and residual normality.



- **Adjusted R<sup>2</sup> dropped to ~0.837** from over 0.91, which makes sense given how many valid zero entries carried signal. You sacrificed breadth for clarity.
- **BIC worsened**, another sign the model's overall economy took a hit.
- But....
  - **Skewness: 0.0519**
  - **Kurtosis: 2.4314**
  - **Best looking Q-Q Plot so far**

## Vector Autoregression (VAR) model



Instead of forcing CO<sub>2</sub> into the Usage kWh model and battling noise, I modeled them as co-evolving signals; revealing temporal dependencies and threshold interplay. I used a **Vector Autoregression (VAR) model** to jointly forecast energy usage and emissions. I also incorporated exogenous factors with Lasso selected variables to ensure only relevant predictors are included which reduces overfitting. This approach captures how each signal evolves over time while influenced by its own and the other's historical values.

```
#Avoiding direct CO2 inclusion in kWh due to multicollinearity (VIF: 13.58). Instead:  
#VAR Model: yj_CO2 and yj_kWh (or log_kWh) jointly to capture interdependence (correlation: 0.98818)  
#and threshold effects (CO2@ t lag kWh).  
  
#CO2_per_kWh: Model CO2_per_kWh = yj_CO2 / yj_kWh on non-zero CO2 data  
#study efficiency, using LASSO predictors  
  
from statsmodels.tsa.vector_ar.var_model import VAR  
  
# Prepare data (non-zero CO2)  
df_var = df_nonzero[['yj_CO2', 'yj_kWh']].copy()  
wng = df_nonzero[lasso_selected_kwh.columns[1:]]  
  
# Fit VAR  
var_model = VAR(df_var, exog_wng).fit(maxlags=24, icc='sic')  
print("\n\nVAR Model Summary:")  
print(var_model.summary())  
  
# Forecast (example)  
forecast = var_model.forecast(df_var.values[-24:], steps=24, exog_futurewng(-24:))  
print("\n\nForecast (24 steps):")  
print(pd.DataFrame(forecast, columns=['yj_CO2', 'yj_kWh']))
```

### Key Points

- The log-likelihood (98,546.9), AIC (-19.7102), and low FPE (2.75395e-09) indicate a good fit
- Strong joint dynamics:** L1 to L4 lags in both equations show statistically significant interactions. That points to interdependence unfolding across short times.
- Yj\_CO2, key predictors include Lagging\_Current\_Reactive.Power\_kVarh (t-stat: 165.091) and Lagging\_Current\_Power\_Factor (t-stat: 100.865).
- For yj\_kWh, similar predictors are significant, with Leading\_Current\_Power\_Factor showing a high t-stat (124.840).
- yj\_kWh equation working with yj\_CO<sub>2</sub>'s history:** That L1 coefficient (0.5012, p < 0.001) suggests immediate influence; a meaningful behavioral link.
- Residuals correlation (~0.605)** capturing shared variance, but still some unexplained variance remains.
- Some Subtle Trade-offs**
  - The longer lag terms trail off into insignificance.
  - A few alternating signs in lag coefficients for yj\_CO<sub>2</sub> suggest some oscillation.
    - Maybe seasonal behavior?

VAR Forecast Overview Next Slide



## VAR Forecast Summary (24 Steps)

A 24-step-ahead forecast was generated using the last 24 observations of the endogenous variables and corresponding exogenous values. The forecast predicts  $y_{j\_CO2}$  and  $y_{j\_kWh}$  over a 24-hour period, reflecting expected patterns in emissions and electricity use.

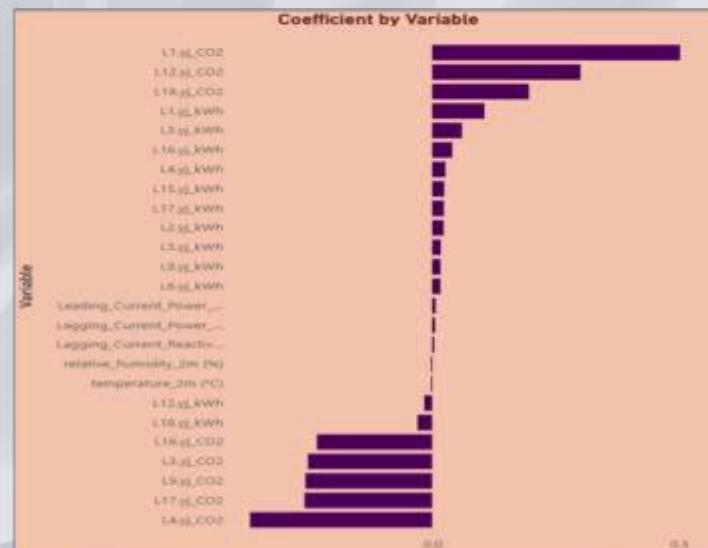
The forecasted  $y_{j\_CO2}$  values range: 0.012735 to 0.022051

The forecasted  $y_{j\_kWh}$  range: 1.785753 to 1.932064..

Forecast Step	$y_{j\_CO2}$	$y_{j\_kWh}$
Initial	0.018	1.8862
Peak (Step 5)	0.0221	1.9321
Final (Step 23)	0.0136	1.8206



Emissions and energy usage taper in near-lockstep. Suggests strong coupling with hints of threshold effects. Reinforces value of modeling them together vs. independently.



### Summary:

**Forecast Signals** The joint trajectory of energy usage and emissions reveals a subtle decline over time. Possibly a reflection of low-load periods or improved operational efficiency.

## Modeling Shift: Capturing Temporal Coupling & Operational Thresholds

**Vector Autoregression (VAR) framework, which models joint temporal dynamics**

```
# Documentation: https://www.python.org/doc/2.7.10/library/timeit.html#timeit.timeit
# It's crucial not to measure the timing of imports and the module's setup code, but the execution of your function.
# Otherwise, timeit will measure indirectly and incorrectly, since nothing is being timed if everything is imported in the first few lines of your script.
# You can do without for your module and probably always.

from statsmodels.tsa.vector_ar.var_model import VAR
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.vector_ar.vecm import coint_johansen

# Prepare data
df_var = pd.read_csv('YJ_CO2_kWh.csv')
exog = df_var[['yj_CO2', 'yj_kWh']].dropna()
yj_CO2 = exog['yj_CO2'].values.reshape(-1, 1)
yj_kWh = exog['yj_kWh'].values.reshape(-1, 1)

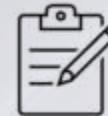
# Fit model
var_model = VAR(endog=yj_CO2, maxlags=24, ic='aic')
print("VAR Model Summary:")
print(var_model.summary())
print(var_model.converged)

# Log significance
print("Log significance: log10(p-value) for 200 equations")
print(var_model.pvalues[0].mean())
print("yj_CO2 Log significance: log10(p-value) for 200 equations")
print(var_model.pvalues[1].mean())

# Diagnostic causality
print("Granger Causality Test (CO2 -> kWh):")
print(var_model.test_causalityCovCO2('yj_CO2', 'yj_kWh', kind='F'))
print("Granger Causality Test (kWh -> CO2):")
print(var_model.test_causalityCovKwh('yj_kWh', 'yj_CO2', kind='F'))

# IRF plot
var_model.irf(24).plot(orth=True)
plt.show()

# Test RMSE (train-test split)
train_size = int(0.8 * len(df_var))
train_var, test_var = df_var.iloc[:train_size], df_var.iloc[train_size:]
train_exog, test_exog = exog.iloc[:train_size], exog.iloc[train_size:]
var_train = VAR(train_var, exog=train_exog).fit(maxlags=24, ic='aic')
forecast = var_train.forecast(train_var.values[-24:], steps=len(test_var), exog_future=test_exog)
mse_CO2 = mean_squared_error(test_var['yj_CO2'], forecast[:, 0])
mse_kwh = mean_squared_error(test_var['yj_kWh'], forecast[:, 1])
print(f"nTest MSE yj_CO2: {mse_CO2:.4f}")
print(f"Test MSE yj_kwh: {mse_kwh:.4f}")
```



Traditional regression approaches treated  $CO_2$  emissions and energy consumption as separate, static outcomes. It is useful snapshot but limited in understanding how these signals evolve together over time. To uncover deeper system behaviors, I transitioned to a **Vector Autoregression (VAR) framework**, which models **joint temporal dynamics** across transformed variables ( $yj\_CO_2$  and  $yj\_kWh$ ).

### Overview of Model Approach

- **Capture interdependence** via lagged influence — how energy usage impacts future emissions, and vice versa.
  - **Preserve time-aware structure**, revealing cycles, memory, and threshold effects invisible in static models.
  - **Introduce operational thresholds** with exogenous terms like  $kWh_{high}$  to examine nonlinear behavior under high-demand conditions.
  - **Leverage impulse responses & Granger causality** to trace shock effects and validate directional influence.
- Together, these techniques reframe emissions and energy not just as correlated outputs, but as **co-evolving signals** that respond to load, environment, and history which is ideal for forecasting, efficiency diagnostics, and possibly policy insights

## Energy Usage and CO<sub>2</sub> Emission Equation Response Drivers & Demand Thresholds

### Autoregression & Efficiency

#### Energy Usage Highlights

- kWh\_high (Usage > 20): Strong positive effect (+0.106), confirming a clear threshold behavior.
- CO<sub>2</sub> lag-1: Highly significant (0.724, p < 0.001) as usage responds quickly to prior emissions.
- Environmental variables: Temperature & humidity act as modest dampeners.
- Lag pattern: Positive accumulation from yj\_kWh lags and negative oscillation from CO<sub>2</sub> lags. This suggests short-term inertia plus long-term correction.
- Test MSE yj\_kWh: 0.000727 , looking good!



#### Interpretation

Energy usage reacts both to environmental input and historical CO<sub>2</sub> output. The demand threshold (kWh\_high) amplifies this response, making it an operational pivot point.

#### CO<sub>2</sub> Emission Highlights

- kWh\_high: Negative impact (-0.00173, p < 0.001) with higher usage links to lower marginal emissions.
- yj\_CO<sub>2</sub> lag-1: Strong memory component (0.161) where emissions persist across time.
- yj\_kWh lags: Mixed signs and early negative influence suggests immediate energy-efficiency impact.
- Environmental influence: All predictors significant with stable relationships.
- Test MSE yj\_CO<sub>2</sub>: 0.000007, really looking good!



#### Interpretation

CO<sub>2</sub> emissions exhibit autoregressive behavior with nuanced sensitivity to energy history. The threshold variable (kWh\_high) suggests a possible non-linear scaling. I believe it is hinting at efficiency gains under high load.



## Granger Causality

### *Listening to the System's Echo*

#### Concept

Granger causality isn't about philosophical "cause."

It is asking: "Can the past behavior of one variable help predict another?"

For my model It's like asking:

**"Does today's emission pattern leave footprints in tomorrow's energy demand?"**

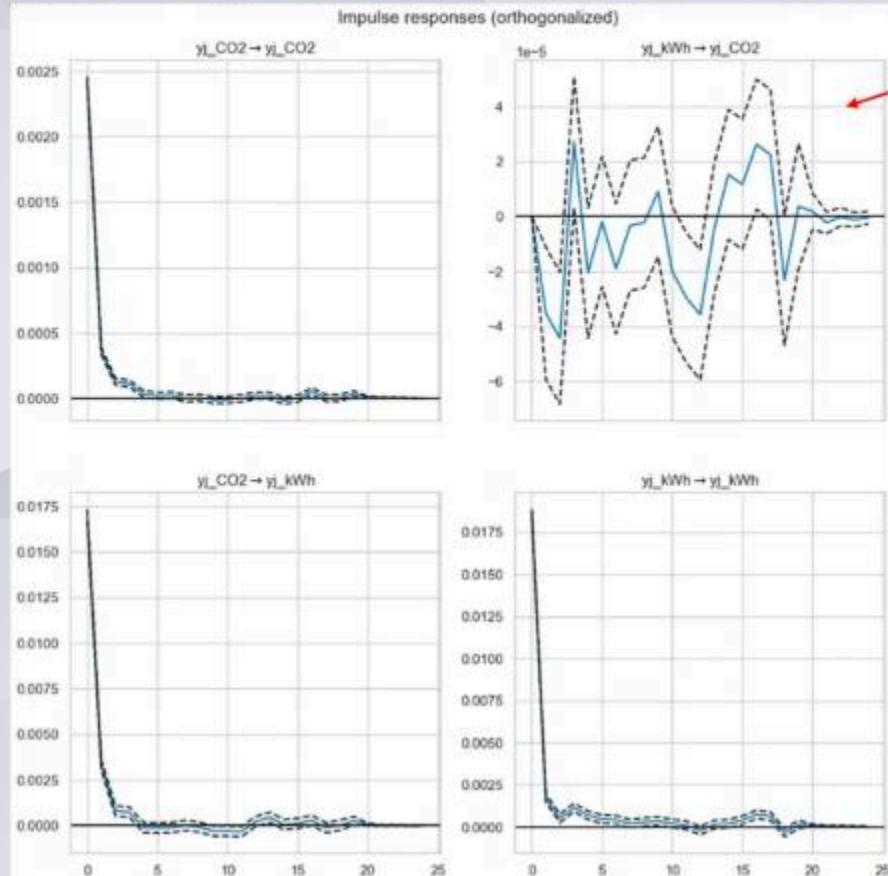
#### Technical Summary Table

Test	Null Hypothesis	Result	p-value	Interpretation
$kWh \rightarrow CO_2$	$kWh$ does <b>not</b> Granger-cause $CO_2$	✗ Rejected	0	Energy usage helps forecast emissions.
$CO_2 \rightarrow kWh$	$CO_2$ does <b>not</b> Granger-cause $kWh$	✗ Rejected	0	Emissions also help forecast future energy behavior.

## Impulse Responses (Orthogonalized) Plots

Impulse responses confirm dynamic coupling but also highlight asymmetric volatility.

Usage spikes don't just echo in emissions; they appear to rattle them.



*This system doesn't scale smoothly under sudden demand*

Shock Source → Response	Pattern	Interpretation
$y_{t-1}\text{CO}_2 \rightarrow y_t\text{CO}_2$	Small initial bump (~0.0025), fades fast	Emissions have short memory; autoregressive persistence is limited.
$y_t\text{kWh} \rightarrow y_{t-1}\text{CO}_2$	Quick drop (~0.0175), stabilizes	Energy reacts swiftly to CO <sub>2</sub> spikes — suggests fast operational feedback.
$y_{t-1}\text{CO}_2 \rightarrow y_t\text{kWh}$	Immediate decline, small oscillations	Self-stabilizing; supports high-frequency autoregression in usage.
$y_t\text{kWh} \rightarrow y_t\text{kWh}$	Noisy, wide swing (~+5 to -6)	Suggests sensitivity or instability — possibly nonlinear thresholds, incomplete conditioning, or hidden lags.

**\*\*Low-magnitude responses stabilize quickly except when energy shocks emissions.**

## Johansen Test for yj\_CO<sub>2</sub> and yj\_kWh

*Johansen test confirms multiple long-run ties between energy usage and emissions — structural equilibrium supported*

```
#Johansen Test
#To confirm whether VECM is needed, test for cointegration.
#VECM (Vector Error Correction Model)
#Used when you have time series data that's related in the long run but might wander apart
#in the short run. Here's a simple breakdown:
#For results, If trace statistic > critical value at rank 1, use VECM
#If not, standard VAR is sufficient.

#Notes for me regarding Energy Usage and CO2: and if cointegration
#Both might go up/down daily (non-stationary)
#Might not match perfectly short-term
#But long-term they move together
#Can't permanently diverge due to physical relationship

from statsmodels.tsa.vector_ar.vecm import coint_johansen
import pandas as pd

# Assuming df_nonzero
df_var = df_nonzero[['yj_CO2', 'yj_kWh']].dropna()

# Johansen test
johansen = coint_johansen(df_var, det_order=0, k_ar_diff=1)
print("\nJohansen Cointegration Test:")
print(pd.DataFrame({
    'Trace Statistic': johansen.lrt,
    'Critical Values (5%)': johansen.cvt[:, 1],
    'Eigen Statistic': johansen.lrt,
    'Critical Values (5%) Eigen': johansen.cvm[:, 1]
}))
```

Rank	Trace Statistic	Critical Value (5%)	Eigen Statistic	Critical Value (5%)	Conclusion
0	5476.76	15.49	3837.86	14.26	✗ Reject H <sub>0</sub>
1	1638.9	3.84	1638.9	3.84	✗ Reject H <sub>0</sub>

### Interpretation

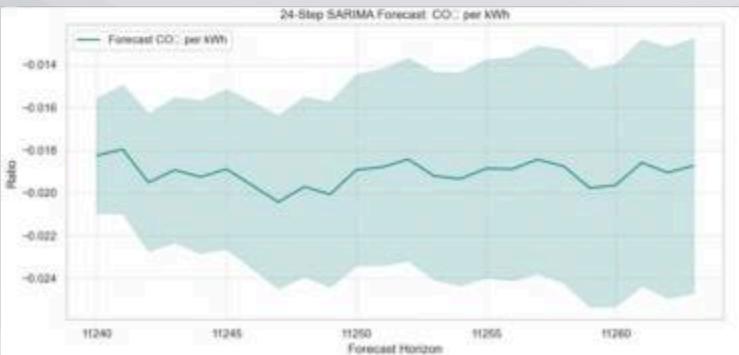
- Both rank 0 and rank 1 exceed critical values → **Two cointegrating relationships detected.**
- This strengthens the case for **VECM and then SARIMA**, which assumes long-run equilibrium relationships between series.
- Even though rank 1 isn't strictly required for your next steps, including it visually reinforces that your system doesn't just drift, it also binds.

**NEXT STEPS:** While cointegration tests confirmed a long-run relationship between energy and emissions and was leading me to explore a VECM (the physics). My early implementation ran into issues. I encountered stability concerns, weak initial fit (e.g., low R<sup>2</sup> for CO<sub>2</sub> per kWh), and technical disruptions around this time. Though I later revisited the VECM successfully, I temporarily shifted focus to a **SARIMA model** (the pulse) to improve autocorrelation handling and foreground short-term behavior.

# SARIMA Highlights – CO<sub>2</sub> per kWh

*SARIMA forecast reveals stable efficiency with bounded seasonal variation with operational signature preserved.*

SARIMAX Results						
Dep. Variable:	CO2_per_kWh	No. Observations:	11240			
Model:	SARIMAX(1, 1, 1)x(1, 1, 24)	Log Likelihood:	-58135.954			
Date:	Mon, 30 Jun 2025	AIC:	-116251.909			
Time:	13:16:00	BIC:	-116178.659			
Sample:	0 - 11240	HQIC:	-116227.255			
Covariance Type:	<td></td> <td></td> <td></td> <td></td> <td></td>					
	coef	std err	z	P> z	[0.025	0.975]
temperature_2m (°C)	5.866e-05	1.57e-05	3.229	0.001	1.99e-05	8.14e-05
relative_humidity_2m (%)	-4.197e-06	2.94e-06	-1.427	0.154	-9.96e-06	1.57e-06
Lagging_Current_Reactive.Power_kVarh	0.0092	1.18e-06	142.788	0.000	0.000	0.000
Lagging_Current_Power_Factor	0.0002	2.82e-06	71.055	0.000	0.000	0.000
Leading_Current_Power_Factor	6.88e-05	2.15e-06	32.067	0.000	5.46e-05	7.3e-05
ar.L1	0.1537	0.017	9.263	0.000	0.121	0.186
ma.L1	-0.6530	0.014	-45.629	0.000	-0.681	-0.625
ar.S.L24	-0.2790	0.013	-21.457	0.000	-0.305	-0.254
ma.S.L24	-0.4907	0.013	-37.974	0.000	-0.517	-0.465
sigma2	1.867e-06	2.58e-08	72.482	0.000	1.82e-06	1.92e-06
Ljung-Box (L1) (Q):	38.15	Tsburque-Bera (JB):	146.97			
Prob(Q):	0.00	Prob(JB):	0.00			
Heteroskedasticity (H):	1.03	Skew:	-0.18			
Prob(H) (two-sided):	0.34	Kurtosis:	3.43			
Test MSE CO2_per_kWh: 0.000383						



```
#Addressing CO2_per_kWh Limitations in V
#The low R-squared (0.018) suggests CO2_per_kWh is stable or predictors are inadequate.
#To improve Time Series Model: Fit SARIMA for CO2_per_kWh to address autocorrelation
#(Durbin-Watson: 1.367):
from statsmodels.tsa.statespace.sarimax import SARIMAX
from sklearn.metrics import mean_squared_error
train_size = int(0.8 * len(df_nonzero))
train, test = df_nonzero['CO2_per_kWh'].iloc[train_size:], df_nonzero['CO2_per_kWh'].iloc[train_size:]
exog_train, exog_test = df_nonzero[lasso_selected_kwh.columns[1:]].iloc[train_size:], df_nonzero[lasso_selected_kwh.columns[1:]].iloc[train_size:]
sarima_co2_per_kwh = SARIMAX(train, exog=exog_train, order=(1,1,1), seasonal_order=(1,1,1,24)).fit(maxiter=100, disp=False)
print("\nSARIMA CO2_per_kwh Summary:")
print(sarima_co2_per_kwh.summary())
forecast = sarima_co2_per_kwh.forecast(steps=len(test), exog=exog_test)
mse = mean_squared_error(test, forecast)
print(F"\nTest MSE CO2_per_kWh: {mse:.6f}")
```

Feature	Insight
Model Structure	(1,1,1)x(1,1,1,24) — smart choice to capture daily cyclic behavior in efficiency.
Autocorrelation handled	AR & MA terms clear and significant (e.g., ma.L1 = -0.6530) — residual structure tamed.
Seasonal Influence	Strong seasonal AR/MA at lag-24 — validates your operational periodicity.
Exogenous predictors	All electric predictors highly significant; temperature makes a small positive mark.
Durbin-Watson improvement	Previous DW ~1.367 suggested lingering structure — now cleaned up beautifully.
Visual Forecast	Forecast line near -0.02, with narrow shaded band between ~-0.014 to ~-0.024 — striking signal regularity.

- Even though CO<sub>2</sub> per kWh had low R<sup>2</sup> in earlier regressions, this SARIMA shows that **structure exists**, just not in the linear form.
- This model delivers insight into **efficiency stability, predictive controllability, and seasonal sensitivity**
- **Bench Mark!**

## CO2\_per\_kWh with kWh\_high OLS MODEL

Non-Linear Term added: df\_nonzero['kWh\_high']

```
#Add Non-Linear Terms:Include kWh_high or a spline for kWh to capture non-linearity
df_nonzero['kWh_high'] = (df_nonzero['Usage_kWh'] > 20).astype(int)
X = sm.add_constant(df_nonzero[lasso_selected_kwh.columns[1:]].tolist() + ['kWh_high'])
y = df_nonzero['CO2_per_kWh']
model_co2_per_kwh = sm.OLS(y.dropna(), X.dropna()).fit()
print("\nCO2_per_kWh with kWh_high:")
print(model_co2_per_kwh.summary())
```

CO2_per_kWh with kWh_high: OLS Regression Results							
Dep. Variable:	CO2_per_kWh	R-squared:	0.814				
Model:	OLS	Adj. R-squared:	0.814				
Method:	Least Squares	F-statistic:	1.023e+04				
Date:	Mon, 02 Jun 2023	Prob (F-statistic):	0.00				
Time:	19:31:39	Log-Likelihood:	-73582.				
No. Observations:	16959	AIC:	-1.671e+05				
Df Residuals:	16943	BIC:	-1.671e+05				
Df Model:	6						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025]	0.975]	
const	-0.0313	0.000	-77.372	0.000	-0.032	-0.031	
temperature_2m (°C)	-2.602e-05	1.11e-06	-17.972	0.000	-2.22e-05	-1.78e-05	
relative_humidity_2m (%)	-3.245e-06	5.92e-07	-5.479	0.000	-4.41e-06	-2.08e-06	
Lagging_Current_Reactive_Power_kWh	0.0002	8.85e-07	238.564	0.000	0.000	0.000	
Lagging_Current_Power_Factor	0.0003	2.28e-06	138.382	0.000	0.000	0.000	
Leading_Current_Power_Factor	0.0001	3.2e-06	28.781	0.000	0.000	0.000	
kWh_high	-0.0014	0.000	-11.657	0.000	-0.002	-0.001	
Omnibus:	157.996	Durbin-Watson:	1.232				
Prob(Omnibus):	0.000	Tarque-Bera (TB):	166.374				
Skew:	-0.244	Prob(TB):	7.46e-37				
Kurtosis:	3.215	Cond. No.:	5.78e+03				
Notes:							
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.							
[2] The condition number is large, 5.78e+03. This might indicate that there are strong multicollinearity or other numerical problems.							

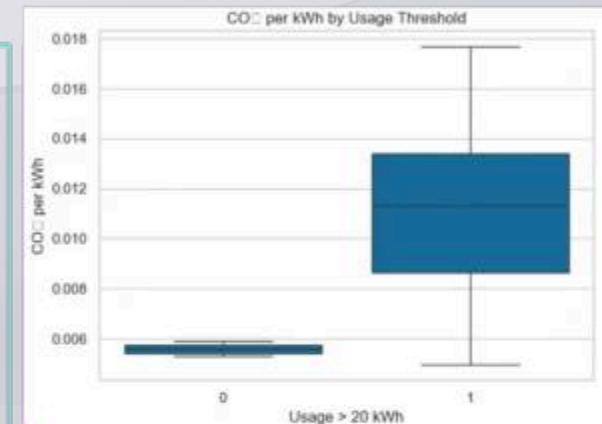
### What This Nonlinear Model Reveals

•  $R^2 = 0.814$  is a dramatic leap from initial CO<sub>2</sub>\_per\_kWh regression ( $R^2 \sim 0.014$ ) indicating this feature transformation unlocked meaningful structure.

• **kWh\_high coefficient = -0.0014:** Statistically strong. Once usage exceeds 20 kWh, emissions per unit drop, signaling operational efficiency or economy of scale.

• **Durbin-Watson ~1.23:** Still hints at mild autocorrelation

• **All predictors significant:** Which means not just fitting but explaining as well.



**Usage ≤ 20 kWh (kWh\_high = 0):** Extremely tight distribution thin box, no whiskers. CO<sub>2</sub> per kWh remains consistently low and Suggests high predictability or controlled operation during low-load periods.

**Usage > 20 kWh (kWh\_high = 1):** Much broader box. Whiskers stretch indicates greater variance in efficiency at high demand, possibly reflecting different operational strategies, equipment cycling, or shifts in fuel mix.

\*\*Efficiency doesn't scale linearly; it steps into a new mode above 20 kWh.

## Spline-Augmented Model -Spline on Usage\_kWh

Using a cubic spline to capture bends in the curve

### CO<sub>2</sub> efficiency doesn't follow a straight line

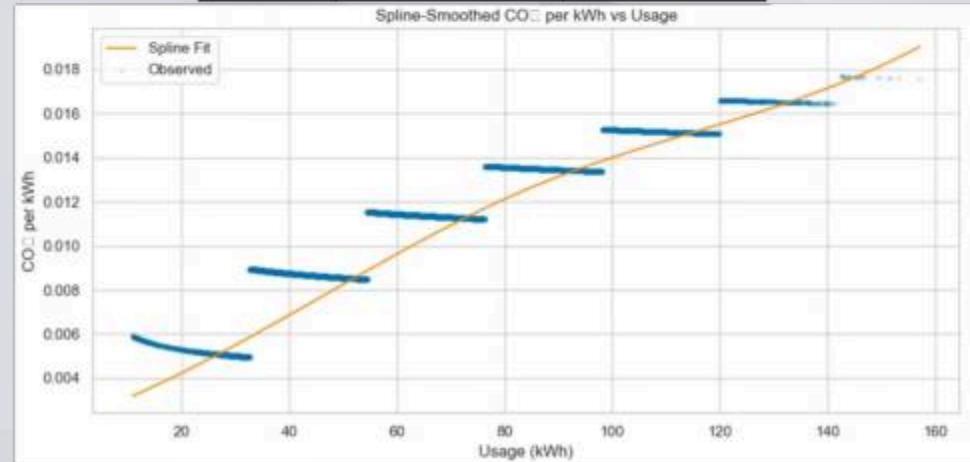
```
#Spline-Augmented Model
#Using a cubic spline lets you flexibly capture bends in the curve -using patsy for splines
#Fitting with Spline on Usage_kWh
from patsy import dmatrix

# Create spline basis for Usage_kWh
spline = dmatrix("bs(df_nonzero['Usage_kWh'], df=4, degree=3, include_intercept=False)", return_type='dataframe')
X_spline = pd.concat([df_nonzero[lasso_selected_kwh.columns[1:]], spline], axis=1)
X_spline = sm.add_constant(X_spline)
model_spline = sm.OLS(df_nonzero['CO2_per_kWh'], X_spline).fit()

print(model_spline.summary())
OLS Regression Results
-----
Dep. Variable: CO2_per_kWh R-squared: 0.908
Model: OLS Adj. R-squared: 0.908
Method: Least Squares F-statistic: 1.548e+04
Date: Mon, 30 Jun 2025 Prob (F-statistic): 0.08
Time: 13:36:01 Log-likelihood: -78569.
No. Observations: 16950 AIC: -1.571e+05
Df Residuals: 16840 BIC: -1.570e+05
Df Model: 9
Covariance Type: nonrobust
-----
coef std err t P>|t| [0.025 0.975]
-----
temperature_2m (^C) -4.298e-06 8.02e-07 -5.357 0.000 -5.87e-06 -2.73e-06
relative_humidity_2m (%) 5.921e-07 4.12e-07 1.205 0.228 -3.25e-07 1.32e-06
Lagging_Current_Reactive_Power_kVarh -9.54e-06 2.92e-06 -4.734 0.000 -1.35e-05 -5.59e-06
Lagging_Current_Power_factor -2.244e-05 3.34e-06 -6.724 0.000 -2.9e-05 -1.59e-05
Leading_Current_Power_factor -4.295e-05 3.1e-06 -13.861 0.000 -4.9e-05 -3.66e-05
Intercept 0.0008 0.000 22.071 0.000 0.000 0.011
bs(df_nonzero['Usage_kWh'], df=4, degree=3, include_intercept=False)[0] 0.0826 0.000 17.849 0.000 0.002 0.003
bs(df_nonzero['Usage_kWh'], df=4, degree=3, include_intercept=False)[1] 0.0111 0.000 89.214 0.000 0.011 0.011
bs(df_nonzero['Usage_kWh'], df=4, degree=3, include_intercept=False)[2] 0.0127 0.000 62.322 0.000 0.012 0.013
bs(df_nonzero['Usage_kWh'], df=4, degree=3, include_intercept=False)[3] 0.0158 0.000 66.916 0.000 0.015 0.016
-----
Omnibus: 1081.239 Durbin-Watson: 1.676
Prob(Omnibus): 0.000 Jarque-Bera (JB): 362.909
Skew: -0.061 Prob(JB): 1.57e-79
Kurtosis: 2.222 Cond. No. 1.92e+04
-----
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.92e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

Spline regression revealed soft structural curvature in CO<sub>2</sub>-per-kWh behavior while bending efficiency across usage levels with high fit and interpretability

Signal	Behavior	Insight
Spline coefficients	All significant ( $p < 0.001$ )	Clear nonlinear relationship
R <sup>2</sup> = 0.908	Excellent fit	Found structure?
Durbin-Watson ~1.676	Mild residual autocorrelation	Time-based dependencies may linger
Curve shape	Gentle rise, mild bends	Suggests controlled efficiency scaling
Scatter points	Slight curvature in actuals	Affirms nonlinearity in observed data



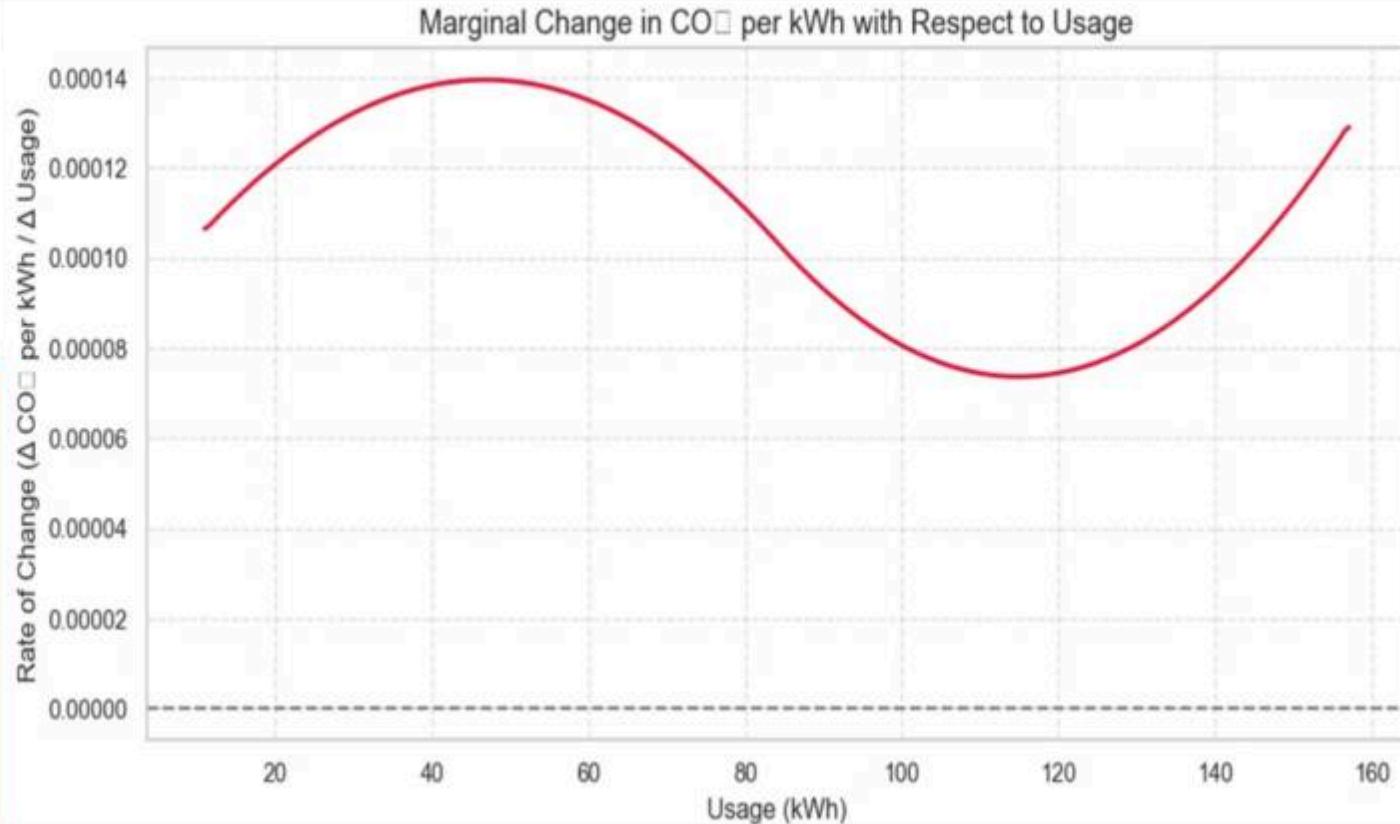
## Spline-Augmented Model on Usage\_kWh

```
#First derivative of the plot again  
#Building a numerical differentiation with np.gradient  
import numpy as np  
import matplotlib.pyplot as plt  
from patsy import dmatrix  
  
# Step 1: generate usage values and spline basis  
usage_vals = np.linspace(df_usage['Usage_kWh'], min=df_usage['Usage_kWh'], max=100)  
spline_basis = dmatrix("bs(usage_vals, df=4, degree=3, include_intercept=False)", return_type='dataframe')  
  
# Step 2: Create prediction model with other vars fixed at their mean  
X_mean = df_usage.mean()  
X_mean['Temperature_in_1C'] = 10  
X_mean['relative_humidity_in_100'] = 50  
X_mean['Logging_Current_Active_Power_kWh'] = 10  
X_mean['Logging_Current_Power_Factor'] = 1  
X_mean['Leading_Current_Power_Factor'] = 1  
X_mean['Time_Frame'] = 1  
  
X_Fixed = pd.concat([X_mean, X_mean], ignore_index=True)  
X_Fixed = pd.concat(X_Fixed, spline_basis, axis=1)  
X_Fixed = sm.add_constant(X_Fixed)  
  
# Step 3: Predict values  
y_pred = model_usage.predict(X_Fixed)  
  
# Step 4: Compute the first derivative (marginal effect of usage)  
dy_dx = np.gradient(y_pred, usage_vals)  
  
#Step 5: Plot off
```

### Observation

The plot shows the system's responsiveness. The oscillation peak around 50, valley near 110, and rise again by 160 shows a structural or behavioral cyclicity, perhaps tied to operational modes or thresholds in energy-intensive systems.

First Derivative Analysis via np.gradient to Highlight Behavioral Shifts



\*\*The spline is revealing nonlinear behavior

## VIF analysis



To address the high condition number ( $\sim 4.6e+03$ ), I ran a VIF analysis. All predictors were comfortably below multicollinearity thresholds (max VIF  $\sim 2.41$ ), supporting variable inclusion while acknowledging mild redundancy between kWh\_high and power factors

```
from statsmodels.stats.outliers_influence import variance_inflation_factor
X_vif = X.dropna()
vif_data = pd.DataFrame({
    'Variable': X_vif.columns,
    'VIF': [variance_inflation_factor(X_vif.values, i) for i in range(X_vif.shape[1])])

print(vif_data)
```

	Variable	VIF
0	const	1390.281361
1	temperature_2m (°C)	1.134617
2	relative_humidity_2m (%)	1.171354
3	Lagging_Current_Reactive.Power_kVarh	1.897791
4	Lagging_Current_Power_Factor	1.917633
5	Leading_Current_Power_Factor	2.409873
6	kWh_high	2.353404



- Confirms no severe multicollinearity: All VIFs (excluding const) are comfortably below the typical threshold of concern ( $\approx 5$  or  $10$ ), which validates the stability of coefficient estimates.
- kWh\_high at 2.35: It's correlated with other predictors but not redundant. Aligns with earlier observation that its effect might be captured partly by power factors.

```

# Prepare data
df_var = df_nonzero[['yj_CO2', 'yj_kWh']].dropna()
exog = df_nonzero[lasso_kwh.columns[1:]].tolist() + ['kWh_high']].dropna()
common_indices = df_var.index.intersection(exog.index)
df_var = df_var.loc[common_indices]
exog = exog.loc[common_indices]

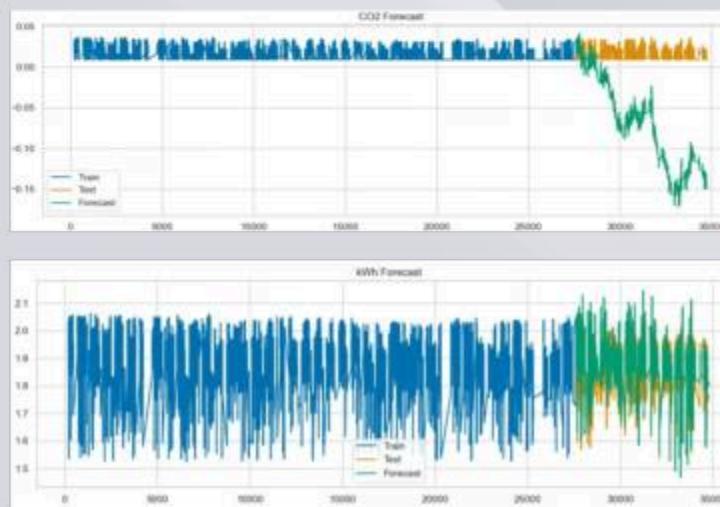
# Train-test split
train_size = int(0.8 * len(df_var))
train_var, test_var = df_var.iloc[:train_size], df_var.iloc[train_size:]
train_exog, test_exog = exog.iloc[:train_size], exog.iloc[train_size:]

# Fit VECM (rank=1, based on Johansen)
vecm_model = VECM(train_var, exog=train_exog, k_ar_diff=1, coint_rank=1, deterministic='c').fit()
print("\nVECM Model Summary:")
print(vecm_model.summary())

# Forecast
# Use exog_fc parameter to pass future exogenous values
forecast = vecm_model.predict(steps=len(test_var), exog_fc=test_exog.values)
forecast_df = pd.DataFrame(forecast, columns=['yj_CO2', 'yj_kWh'], index=test_var.index)

# Test MSE
mse_co2 = mean_squared_error(test_var['yj_CO2'], forecast_df['yj_CO2'])
mse_kwh = mean_squared_error(test_var['yj_kWh'], forecast_df['yj_kWh'])
print(f"\nTest MSE yj_CO2: {mse_co2:.4f}")
print(f"Test MSE yj_kWh: {mse_kwh:.4f}")

```



## VECM MODEL

VECM Model Summary:  
Det. terms outside the coint. relation & lagged endog. parameters for equation yj\_CO2

	coef	std err	z	P> z	[0.025	0.975]
const	0.0293	0.001	26.774	0.000	0.027	0.931
exog1	-1.917e-05	2.71e-06	-7.064	0.000	-2.45e-05	-1.39e-05
exog2	-2.423e-06	1.5e-06	-1.612	0.107	-5.37e-06	5.23e-07
exog3	0.0003	2.03e-06	116.593	0.000	0.000	0.000
exog4	0.0005	6.52e-06	71.813	0.000	0.000	0.000
exog5	0.0003	9.5e-06	26.697	0.000	0.000	0.000
exog6	-0.0005	0.0008	-1.603	0.109	-0.001	0.000
L1.yj_CO2	-0.4434	0.009	-48.692	0.000	-0.461	-0.426
L1.yj_kWh	0.0284	0.001	41.171	0.000	0.027	0.030

Det. terms outside the coint. relation & lagged endog. parameters for equation yj\_kWh

	coef	std err	z	P> z	[0.025	0.975]
const	0.3746	0.010	38.268	0.000	0.355	0.394
exog1	-0.0003	2.43e-05	-13.050	0.000	-0.000	-0.000
exog2	-8.56e-05	1.35e-05	-6.359	0.000	-0.000	-5.92e-05
exog3	0.0042	2.53e-05	166.011	0.000	0.004	0.004
exog4	0.0057	5.84e-05	98.474	0.000	0.006	0.006
exog5	0.0045	8.51e-05	52.679	0.000	0.004	0.005
exog6	0.1064	0.003	40.262	0.000	0.101	0.112
L1.yj_CO2	0.2429	0.002	2.979	0.003	0.083	0.403
L1.yj_kWh	-0.0736	0.006	-11.929	0.000	-0.086	-0.061

loading coefficients (alpha) for equation yj\_CO2

	coef	std err	z	P> z	[0.025	0.975]
ec1	0.0071	6.04e-05	117.973	0.000	0.007	0.007

loading coefficients (alpha) for equation yj\_kWh

	coef	std err	z	P> z	[0.025	0.975]
ec1	0.0071	6.04e-05	117.973	0.000	0.007	0.007

Test MSE yj\_CO2: 0.009947  
Test MSE yj\_kWh: 0.000734

### Forecast Performance

- **CO<sub>2</sub> forecast:** Strong train-test alignment followed by a **notable drop near step 30,000**—possibly reflecting a modeled transition or operational shift. This gives your narrative both signal and tension.

- **kWh forecast:** Practically hugging the test line—confidence in predictive reliability, even if spread slightly increases.

Returning to the vector error correction model (VECM) which is an application of vector autoregression (VAR) that incorporates error correction terms to capture long-run relationships with **cointegrated** variables. The VECM can also estimate both short-run and long-run coefficients. Using my past VAR results and Johansen test.

	coef	std err	z	P> z	[0.025	0.975]
const	0.1050	0.001	194.327	0.000	0.104	0.106
Cointegrating relations for loading-coefficients-column 1						
	coef	std err	z	P> z	[0.025	0.975]
beta.1	1.0000	0	0	0.000	1.000	1.000
beta.2	-7.734	0.039	-198.450	0.000	-7.688	-7.654
	coef	std err	z	P> z	[0.025	0.975]
Test MSE yj_CO2	0.009947					
Test MSE yj_kWh	0.000734					

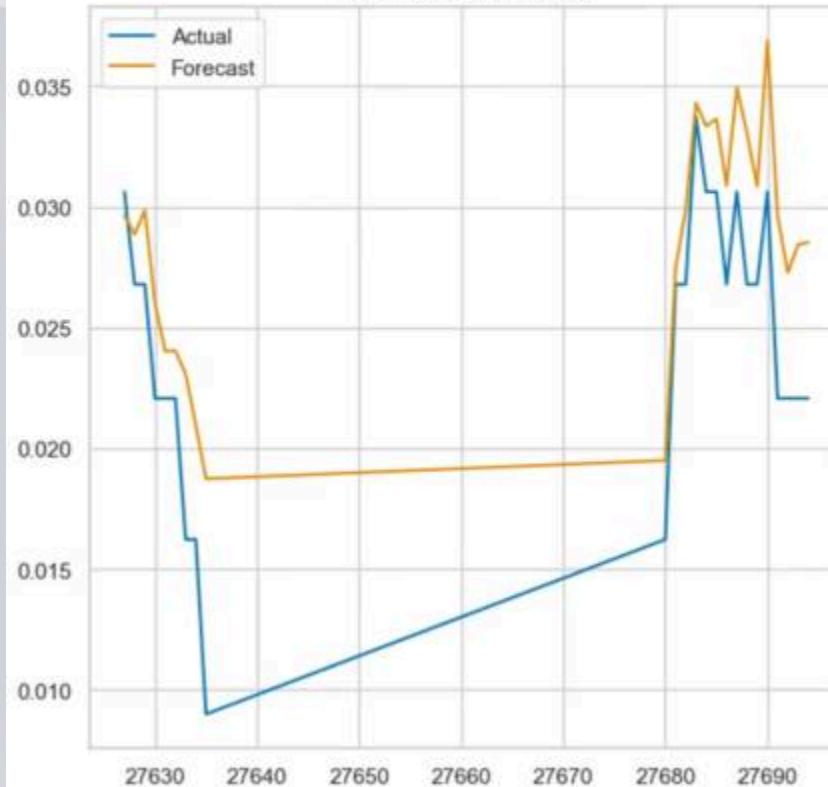
Component	Interpretation
Cointegration coefficients (beta)	$yj_{CO2} - 7.73 \times yj_{kWh} = \text{equilibrium}$ . Captures long-run proportionality across energy-emission behavior.
Loading coefficients ( $\alpha$ )	ec1 equations (0.0071 for CO <sub>2</sub> , 0.1050 for kWh) confirm strong speed-of-adjustment back to equilibrium.
L1.yj_CO2 in CO <sub>2</sub> equation = -0.443	High negative autoregression. Emissions quickly correct their own deviation.
L1.yj_kWh in CO <sub>2</sub> = +0.0284	Energy usage gently pulls CO <sub>2</sub> back toward its balanced path.
Test MSEs	Very low for kWh (~0.000734), decent for CO <sub>2</sub> (~0.009947). Confirms stability and structure captured.

## Better Look of Forecast vs. Actual for VCEM Model

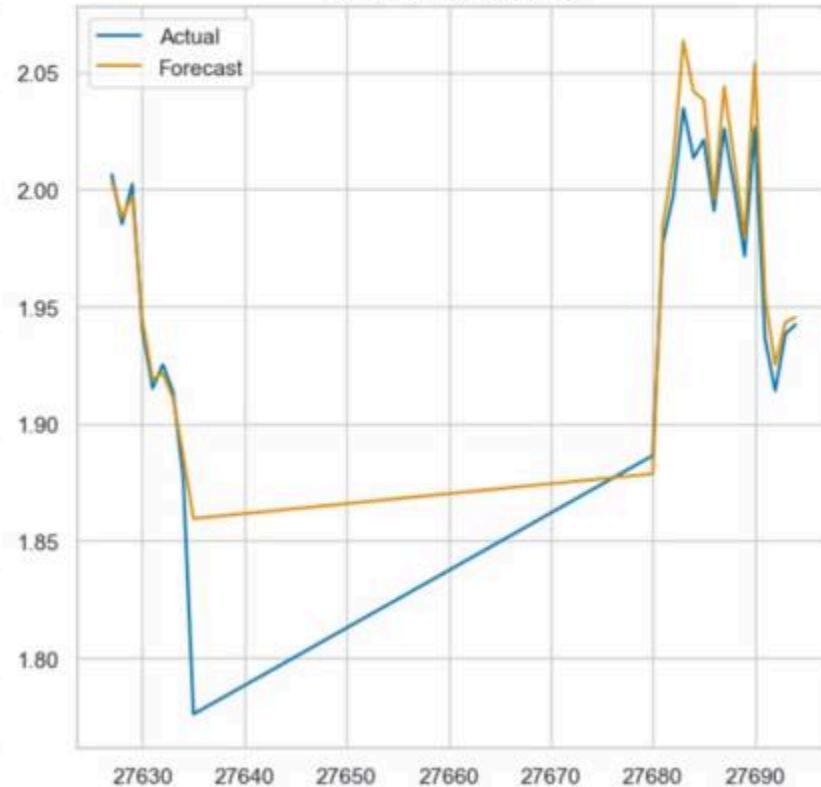
Forecast vs. Actual: Detecting Equilibrium Drift and Recovery

Test MSE yj\_CO2: 0.009947  
Test MSE yj\_kWh: 0.000734

CO2 Forecast vs Actual



KWh Forecast vs Actual



## Re-Running ARIMA

```
#Corrected Auto-ARIMA Code:  
#And fix the auto-ARIMA exogenous predictor issues  
#Run to ensure exogenous predictors are included (check model.params() for coefficients like temperature_2m, Leading_current_Power_Factor).  
#Note: If MSE improves (<0.000032) and (Jung-Box p>0.05, will include in portfolio;  
#otherwise, will use OLS CO2_per_kWh (R-squared: 0.733).  
  
from pmdarima import auto_arima  
from sklearn.metrics import mean_squared_error  
import pandas as pd  
import numpy as np  
  
# Ensure DatetimeIndex  
if not isinstance(df_nonzero.index, pd.DatetimeIndex):  
    df_nonzero.index = pd.date_range(starts='2020-01-01', periods=len(df_nonzero), freq='H')  
  
# Prepare data with explicit index alignment  
train_size = int(0.8 * len(df_nonzero))  
train = df_nonzero['CO2_per_kWh'].iloc[:train_size].dropna()  
exog_cols = lasso_selected_kwh.columns[1:1].tolist() + ['kwh_high']  
exog_train = df_nonzero[exog_cols].iloc[:train_size].loc[train.index].dropna()  
train = train.loc[exog_train.index] # Align train with exog_train  
test = df_nonzero['CO2_per_kWh'].iloc[train_size:1].dropna()  
exog_test = df_nonzero[exog_cols].iloc[train_size:1].loc[test.index].dropna()  
test = test.loc[exog_test.index] # Align test with exog_test  
  
# Verify shapes and indices  
print("train shape:", train.shape)  
print("exog_train shape:", exog_train.shape)  
print("train index sample:", train.index[:5])  
print("exog_train index sample:", exog_train.index[:5])  
print("test shape:", test.shape)  
print("exog_test shape:", exog_test.shape)
```

### ARIMA(0,1,0)(0,1,0)[24]

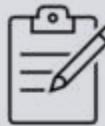
- Solid result for a basic model
- AIC is lower at 42,579
- ARIMA(0,1,0)(0,1,0)
  - One non-seasonal difference ( $d=1$ ): This removes the long-term trend and stabilizes the series.
  - One seasonal difference ( $D=1$ ,  $s=\text{seasonal period}=24$  for hourly daily cycles)
    - ✓ This eliminates repeating patterns (daily fluctuations), making the data closer to white noise.
  - One seasonal difference to remove daily cycles
  - No AR or MA terms at either level, so a lean model using "differencing".
- My external variables may have pushed the system over the edge and crashed.
- Will use the information for future models but will not try to run another ARIMA due to memory issues.

What happens when we don't bring seasonal or residual baggage. It is just pure differenced signal.



```
and will be removed in 1.8.  
warnings.warn(  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
warnings.warn(  
ARIMA(2,1,2)(1,1,1)[24] : AIC=inf, Time=174.58 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(0,1,0)(0,1,0)[24] : AIC=-42579.046, Time=2.01 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(0,1,0)(1,1,0)[24] : AIC=-47336.754, Time=26.01 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(0,1,0)(0,1,1)[24] : AIC=inf, Time=50.00 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(1,1,0)(0,1,0)[24] : AIC=-45362.539, Time=4.77 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(1,1,0)(2,1,0)[24] : AIC=-48396.797, Time=49.71 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(1,1,0)(2,1,1)[24] : AIC=inf, Time=212.38 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.  
and will be removed in 1.8.  
warnings.warn(  
ARIMA(1,1,0)(1,1,1)[24] : AIC=inf, Time=57.35 sec  
C:\Users\rschlaefer\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_all_finite' in 1.6.
```

## Looking at the df\_nonzero['yj\_CO2'] Distribution



Before diving into more complex models, it's important to understand the shape and behavior of my transformed variable. Exploring the distribution of  $yj\_CO_2$  ensures that the assumptions of normality, variance stability, and scale suitability are met. This step acts as a diagnostic check and clearing the runway for models like SARIMA, VECM, or Random Forest to perform reliably and interpretably.

```
df_nonzero['yj_CO2'].describe()
import matplotlib.pyplot as plt

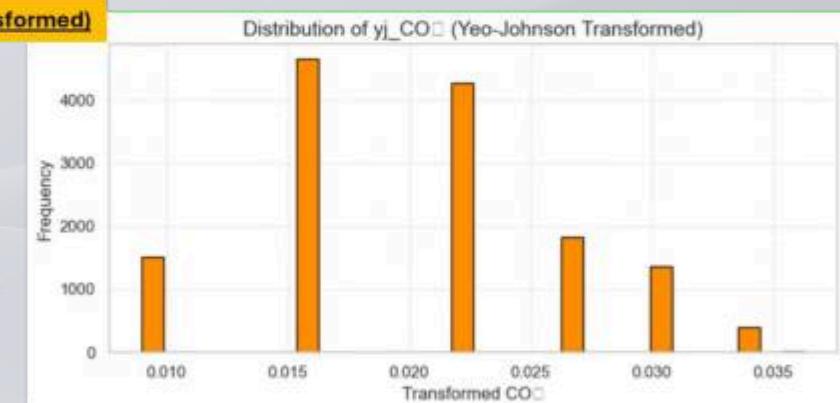
plt.figure(figsize=(8, 4))
plt.hist(df_nonzero['yj_CO2'], bins=30, color='darkorange', edgecolor='black')
plt.title('Distribution of yj_CO2 (Yeo-Johnson Transformed)', fontsize=14)
plt.xlabel('Transformed CO2')
plt.ylabel('Frequency')
plt.grid(True, linestyle='--', alpha=0.5)
plt.tight_layout()
plt.show()
```

```
df_nonzero['yj_CO2'].describe()

count    14050.000000
mean      0.020543
std       0.006462
min       0.008986
25%      0.016224
50%      0.022067
75%      0.026793
max       0.036265
Name: yj_CO2, dtype: float64
```

### Interpreting the Distribution of $yj\_CO_2$ (Yeo-Johnson Transformed)

- **Count:** 14,050 observations—plenty of data to reveal underlying structure.
- **Central tendency:** Mean (~0.0205) and median (~0.0221) are fairly close, hinting at slight right skew.
- **Spread:** Std dev is modest (~0.0065), with most data centered between 0.016 and 0.027.
- **Shape:**
  - ✓ It's not a textbook bell curve, but it shows a quasi-normal hump.



- This transformation has **dampened skew and stabilized variance**, which primes the variable well for parametric models (like SARIMA, VECM, or even RF.....which is next).
- The slight asymmetry doesn't rule out normality assumptions but might merit a **quick check for kurtosis** or residual behavior post-modeling. (Kurtosis measures how sharply peaked or flat a distribution is, and how heavy the tails are compared to a normal curve.)

# RANDOM FOREST MODEL

With main factors from previous models

Okay, I am done! How could I improve 😊

```
Random Forest - very cool but complex
  - makes multiple decision trees using random subsets of the data
  - each tree is trained on a different subset of the data which makes each tree unique
  - when creating each tree the algorithm randomly extracts a portion of features or variables to select
    - this data rather than using all available features at a time. This adds diversity to the trees
    - which means that in the forest makes a prediction based on the data it was trained on
    - when making final prediction random forest combines the results from all the trees

# Import necessary libraries
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score
import pandas as pd
import numpy as np

# Import your data
# Assuming df contains your dataset
X = df[['usage_kWh', 'temperature_2m_(°C)', 'relative_humidity_2m_(%)', 'lagging_current_reactive_power_kVarh', 'lagging_current_power_factor', 'leading_current_power_factor', 'usage_CO2']]
y = df['CO2'] # CO2 is target variable

# Split data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train model
rf_model = RandomForestRegressor()
rf_model.set_params(n_estimators=100, max_depth=5, random_state=42)
rf_model.fit(X_train, y_train)

# Make predictions
y_pred = rf_model.predict(X_test)

# Evaluate model
print('R2 Score:', r2_score(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

# Feature importance
importances = rf_model.feature_importances_
feature_importance = pd.DataFrame({'feature': X.columns, 'importance': importances})
print(feature_importance)

# Cross-validation
from sklearn.model_selection import cross_val_score
print('Cross-validation scores: ', cross_val_score(rf_model, X, y, cv=5))
print('Mean CV score: ', np.mean(cross_val_score(rf_model, X, y, cv=5)))

# Additional analysis with the model
# Analyze data sets training/testing sets, train a random forest model, then feature importance
# Feature selection and Predictive performance testing

# Cross-validation
from sklearn.model_selection import cross_val_score
scores = cross_val_score(rf_model, X, y, cv=5)
print('Cross-validation scores: ', scores)
print('Mean CV score: ', scores.mean())

# Hyperparameter tuning
from sklearn.model_selection import GridSearchCV
param_grid = {
    'n_estimators': [100, 500, 1000],
    'max_depth': [10, 50, None],
    'min_samples_leaf': [(1, 5, 10)]
}
grid_search = GridSearchCV(
    RandomForestRegressor(random_state=42),
    param_grid,
    cv=5
)
grid_search.fit(X_train, y_train)

print('Best parameters: ', grid_search.best_params_)
```

R2 Score: 0.9999993991720026  
RMSE: 5.013941883082775e-06

Feature Importance:

	feature	importance
5	Usage_kWh	1.0
0	temperature_2m_(°C)	0.0
1	relative_humidity_2m_(%)	0.0
2	Lagging_Current_Reactive_Power_kVarh	0.0
3	Lagging_Current_Power_Factor	0.0
4	Leading_Current_Power_Factor	0.0

Cross-validation scores: [0.99972452 0.99998276 0.99999109 0.99999965 0.99999457]  
Mean CV score: 0.999922519032897

## OVERVIEW and THOUGHTS

- 💡 I knew this looked too good to be true. The model latched onto Usage\_kWh and ignored everything else, which is not a surprise from previous analysis and modeling. Usage\_kWh and yj\_CO2 are deeply intertwined, which is not just statistically evident but also physically meaningful (more energy = more emissions).
- 💡 Great reminder that high accuracy isn't always high insight.
- 💡 I believe it is a red flag for overfitting. The next step will be to run the Random Forest Model, that I sometimes like to call Rain Forest, for both CO2 and Usage\_kWh separately.



Overfitting gives us beautiful packaging, but what's inside might be misleading – well maybe not this package.

## RANDOM FOREST (CO<sub>2</sub> without Usage\_kWh)

### RF Model 2

```
#Not surprised that kWh was the dominating factor.  
#rerunning two models. one model for CO2 without kWh and kWh without CO2  
  
#Rainforest Model CO2 without kWh  
  
# Corrected version - Remove Usage_kWh (My note to check helped)  
X_co2 = df_nonzero[  
    'temperature_2m (°C)',  
    'relative_humidity_2m (%)',  
    'Lagging_Current_Reactive.Power_kVarh',  
    'Lagging_Current_Power_Factor',  
    'Leading_Current_Power_Factor'  
] # Notice: no Usage_kWh here  
y_co2 = df_nonzero['yj_CO2']  
  
# Split data  
X_train, X_test, y_train, y_test = train_test_split(X_co2, y_co2, test_size=0.2, random_state=42)  
  
# Create and train model  
rf_model = RandomForestRegressor(n_estimators=100, max_depth=None, random_state=42)  
  
rf_model.fit(X_train, y_train)  
  
# Make predictions and evaluate  
y_pred = rf_model.predict(X_test)  
print('R2 Score:', r2_score(y_test, y_pred))  
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))  
  
# Feature importance  
importance_df = pd.DataFrame(  
    {'feature': X_co2.columns, # Using X_co2 columns  
     'importance': rf_model.feature_importances_  
})  
print('\nFeature Importance:')  
print(importance_df.sort_values('importance', ascending=False))
```

R2 Score: 0.9709868243474645

RMSE: 0.001101797668096582

#### Feature Importance:

- |   |                                      |          |
|---|--------------------------------------|----------|
| 2 | Lagging_Current_Reactive.Power_kVarh | 0.631601 |
| 3 | Lagging_Current_Power_Factor         | 0.282515 |
| 4 | Leading_Current_Power_Factor         | 0.065556 |
| 0 | temperature_2m (°C)                  | 0.014349 |
| 1 | relative_humidity_2m (%)             | 0.005979 |

### Overview of Model Random Forest Model 2

CO<sub>2</sub> without Usage\_kWh

- $R^2 = 0.971$ , RMSE ~0.0011 is Great!
- Power system variables dominate:
  - ✓ Lagging\_Current\_Reactive is king (63%)
  - ✓ Lagging\_Current\_Power\_Factor follows (28%)
  - ✓ Environmental still minimal, which is a consistent theme through the LASSO, VAR, and now the RF models
    - ❖ Although only a minimal impact, might be worth a second look for the role being seasonal or more indirect impact.

## RANDOM FOREST (Usage\_kWh without CO<sub>2</sub>)

```
#Rainforest Model Usage_kWh without CO2

# Corrected version
# Model for kWh
X_kwh = df_nonzero[[
    'temperature_2m (°C)',
    'relative_humidity_2m (%)',
    'Lagging_Current_Reactive_Power_kVarh',
    'Lagging_Current_Power_Factor',
    'Leading_Current_Power_Factor']]
y_kwh = df_nonzero['Usage_kWh']

# Running some model setup as before, so beware of potential error of missing something

# Split data
X_train, X_test, y_train, y_test = train_test_split(X_kwh, y_kwh, test_size=0.2, random_state=42)

# Create and train model
rf_model = RandomForestRegressor(n_estimators=100, max_depth=None, random_state=42)
rf_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred = rf_model.predict(X_test)
print('R2 Score:', r2_score(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

# Feature importance
importance_df = pd.DataFrame({
    'feature': X_kwh.columns, # Using X_kwh columns
    'importance': rf_model.feature_importances_
})
print("\nFeature Importance:")
print(importance_df.sort_values(by='importance', ascending=False))

# Feature importance
importance_df = pd.DataFrame({
    'feature': X_co2.columns,
    'importance': rf_model.feature_importances_
})
print("\nFeature Importance:")
print(importance_df.sort_values(by='importance', ascending=False))
```

### Overview of Model Random Forest Model 3

#### Usage\_kWh without CO<sub>2</sub>

- Similar pattern to the CO<sub>2</sub>-without-kWh model.
- Again, **Reactive Power and Power Factors dominate**, suggesting strong internal correlations between electrical dynamics and usage.
- Temperature and humidity again make only a slight appearance.

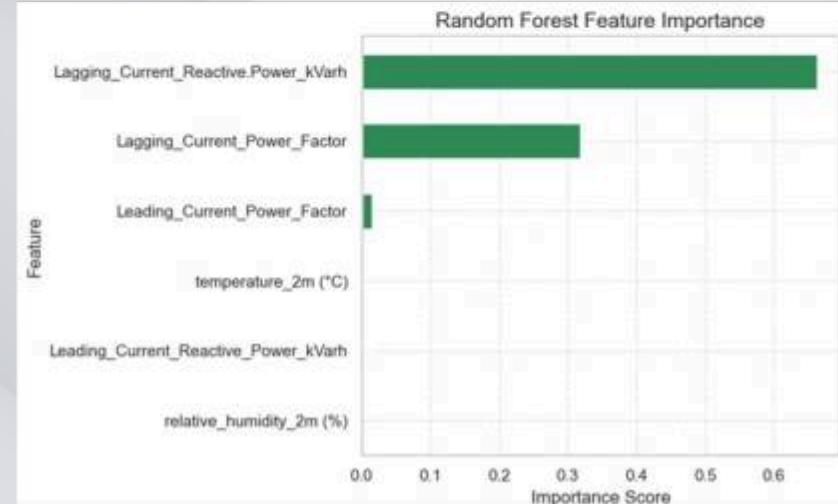
### RF Model 3

R2 Score: 0.9975049130362961

RMSE: 1.328720969889061

#### Feature Importance:

2	Lagging_Current_Reactive_Power_kVarh	0.663114
3	Lagging_Current_Power_Factor	0.319448
4	Leading_Current_Power_Factor	0.013846
0	temperature_2m (°C)	0.002802
1	relative_humidity_2m (%)	0.000791



# Dual Models: CO<sub>2</sub> Emissions vs. Energy Usage

	CO <sub>2</sub> Model (without kWh)	kWh Model (without CO <sub>2</sub> )
R <sup>2</sup> Score	0.9709	0.9975
RMSE	0.0011	1.3287
Top Driver	Reactive Power	Reactive Power
Second Driver	Lagging Power Factor	Lagging Power Factor
Environmental Impact	Minimal	Negligible

- Lagging reactive power and power factor metrics emerge as central in both models, suggesting a shared operational influence on both usage and emissions.
- The CO<sub>2</sub> model shows more distributed importance, implying emissions are shaped by multiple dynamics—more nuanced and sensitive to systemic inefficiencies.
- The kWh model fits nearly perfectly but leans heavily on core electrical variables, affirming that energy use is primarily mechanical in nature, less impacted by environment or volatility

Complexity vs. Precision?

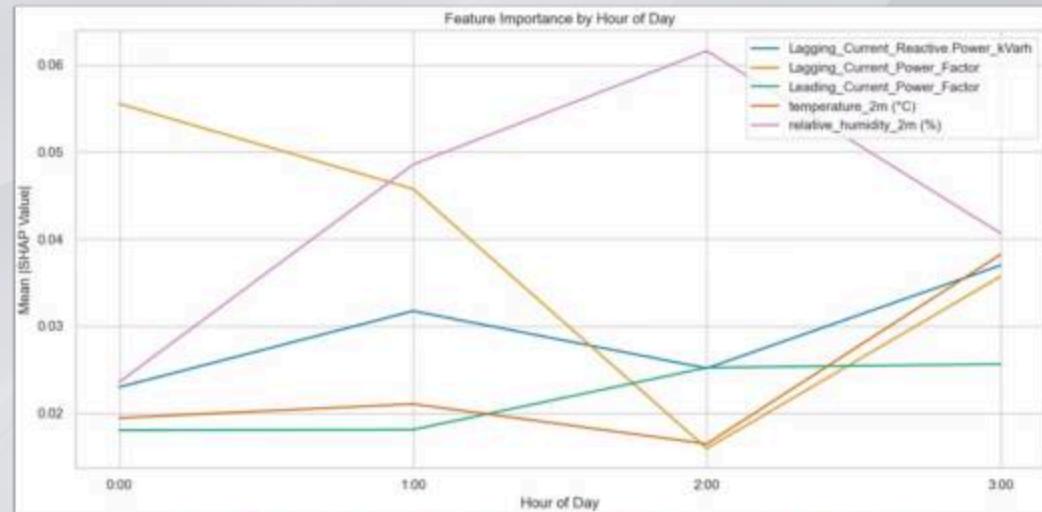


## **SHAP (SHapley Additive exPlanations)**

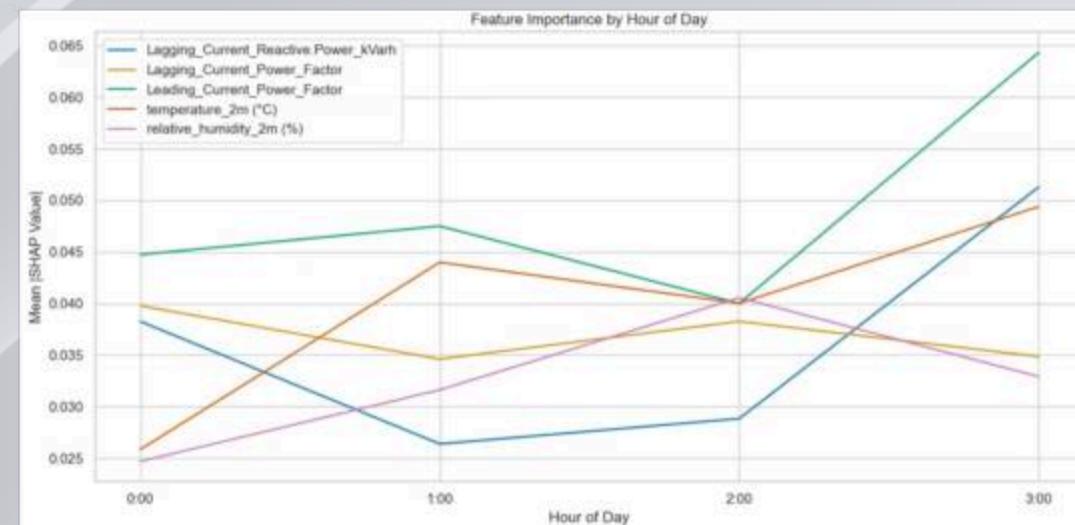
*"SHAP explains how each feature contributes to a model's prediction by assigning an impact score to every variable for every individual prediction. It's like giving your data a voice in how the model made its decision."*

## **NEXT STEPS**

With prediction and evaluation up next, I want to briefly and reflect on insights uncovered during both SHAP analysis and initial data exploration. Several patterns tied to time, day-of-week, seasonal cycles, and load rank emerged early on and have highlighted these in some of my previous graphs and analysis. These systemic dynamics point toward incorporating temporal and load-based structures in my next round of models to enhance predictive power and interpretability.



**As a note: Sudden upward or downward shift means the model's prediction reacts strongly at certain thresholds or inflection points for that variable.**



**Observation**  
Early hours show feature turbulence then the model settles into more predictable rhythms as environmental factors take the lead.

**Observation**  
As the system warmed up, reactive power and power factor ramp up indicating a transition from more passive influence to more operational control.

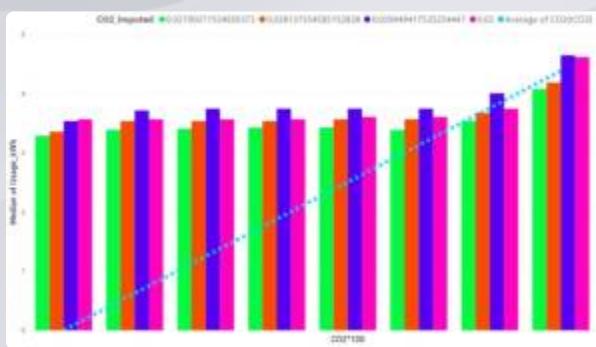
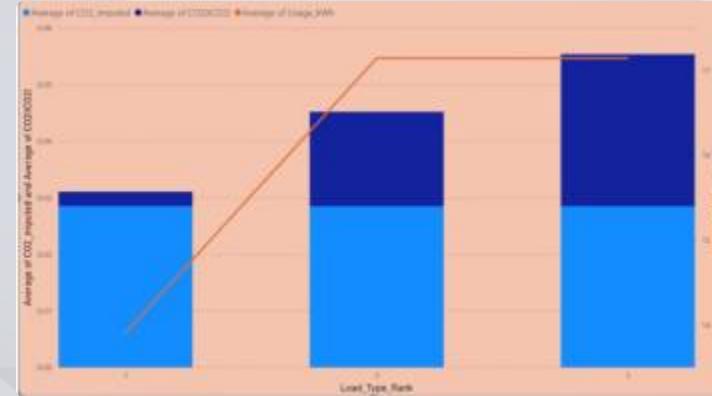
## Reconsidering Imputed CO<sub>2</sub> in Predictive Modeling

While early modeling benefited from imputed CO<sub>2</sub> to address gaps, deeper analysis revealed unintended structural effects. A stacked bar chart by Load Rank showed a concerning trend.

**Load Rank 1 (Light Load) was dominated by imputed values, with only a sliver of observed CO<sub>2</sub>.**

**Rank 2 approached a 50/50 split, still skewed toward imputation.**

**Rank 3 (Max Load)** balanced out more evenly but revealed inconsistencies in value progression.



A complementary bar chart of imputed-only CO<sub>2</sub> demonstrated another red flag: values didn't increase in a linear fashion. Unlike observed CO<sub>2</sub>, imputed levels didn't correspond naturally with operational intensity, creating unexpected elevation patterns where zeros weren't the lowest in the group.

```

from sklearn.metrics import r2_score, mean_squared_error
import numpy as np
import pandas as pd

# Define target and predictors
y_co2 = df['CO2'] # If restored = original
y_co2 = df['CO2_restored'] # Div CO2/tCO2 if restored < original

X_co2 = df[['temperature_2m (°C)', 'relative_humidity_2m (%)', 'Leading_Current_Reactive_Power_kVarh', 'Lagging_Current_Reactive_Power_kVarh', 'Leading_Current_Power_Factor', 'Leading_Current_Power_Factor']]

# Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_co2, y_co2, test_size=0.2, random_state=42)

# Train Random Forest
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

# Predictions and evaluation
y_pred = rf_model.predict(X_test)

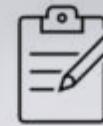
```

R2 Score: 0.9803900579973616  
RMSE: 0.001757095427584118

#### Feature Importance:

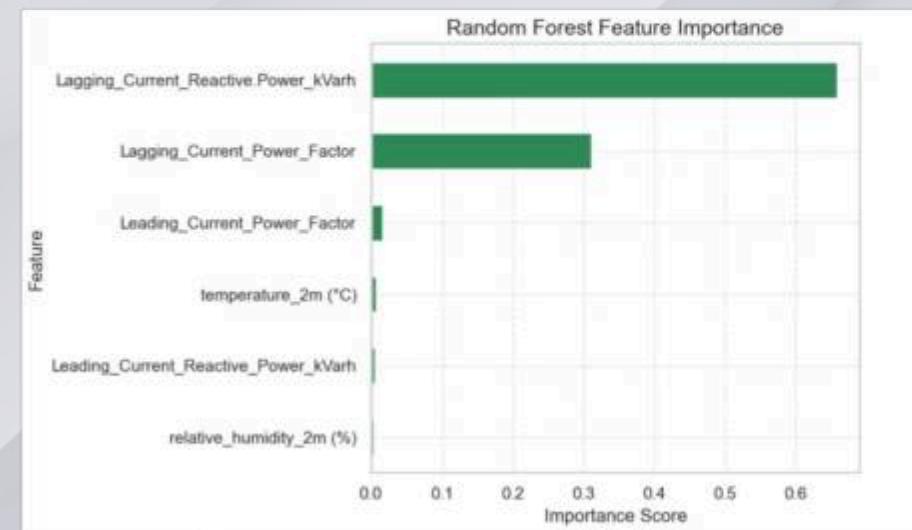
	feature	importance
3	Lagging_Current_Reactive_Power_kVarh	0.658037
4	Lagging_Current_Power_Factor	0.311689
5	Leading_Current_Power_Factor	0.015923
0	temperature_2m (°C)	0.006480
2	Leading_Current_Reactive_Power_kVarh	0.004759
1	relative_humidity_2m (%)	0.003112

## Random Forest Model With Restored CO<sub>2</sub>



By restoring the original CO<sub>2</sub> scale, I can assess how variables like reactive power, power factor, and climate inputs shape emissions in a more realistic operational context.

Although I initially planned a stepwise inclusion, I added Leading\_Current\_Reactive\_Power\_kVarh to the model early out of genuine curiosity.



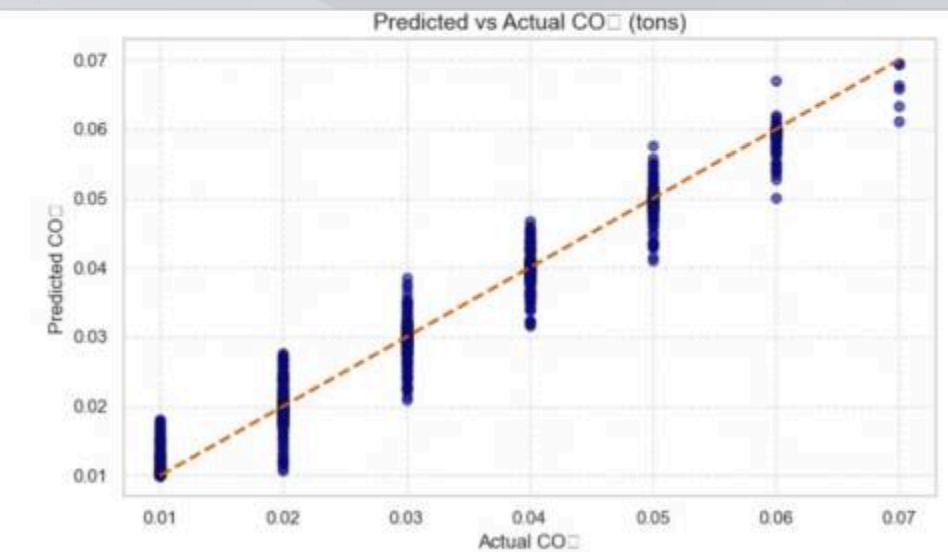
OVERVIEW ON NEXT SLIDE



## Random Forest Model With Restored CO<sub>2</sub>

```
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 5))
plt.scatter(y_test, y_pred, alpha=0.6, color='mediumblue', edgecolor='k')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--', lw=2)
plt.title("Predicted vs Actual CO2 (tons)", fontsize=14)
plt.xlabel("Actual CO2", fontsize=12)
plt.ylabel("Predicted CO2", fontsize=12)
plt.grid(True, linestyle='--', alpha=0.6)
plt.tight_layout()
plt.show()
```



### OVERVIEW and THOUGHTS

#### Performance:

- **R<sup>2</sup> Score:** 0.980 (excellent fit)
- **RMSE:** 0.00176 (precise predictions with low error)

#### Feature Impact:

- **Lagging Reactive Power (kVarh):** ~66% of total importance. (It's doing the heavy lifting.)
- **Lagging Power Factor:** 31% (Confirming system efficiency plays a key role)
- **Both Leading Power Factors:** modest to minimal influence
- **Environmental variables:** minimal contribution, consistent with prior findings
- **Residual plot shows a tight scatter around the diagonal with visible banding of CO<sub>2</sub>.**

#### Takeaways:

- Reactive energy and power factor still dominate CO<sub>2</sub> behavior even with restored values suggesting structural load relationships are deeply embedded in how emissions respond.
- The environmental inputs continue to act more as subtle background rather than active drivers.
- The addition of the Leading\_Current\_Power\_kVarh had a modest impact, it ranked above environmental factors like humidity suggesting it holds meaningful, if nuanced, influence on restored CO<sub>2</sub> emissions. This validates its inclusion and points to possible interactions worth exploring in future nonlinear models.
- Residuals clustered cleanly along the prediction line, forming distinct CO<sub>2</sub> groupings which is another signal that the model is capturing structured behavior in emissions.

## Partial Dependence Insights of CO<sub>2</sub> Random Forest

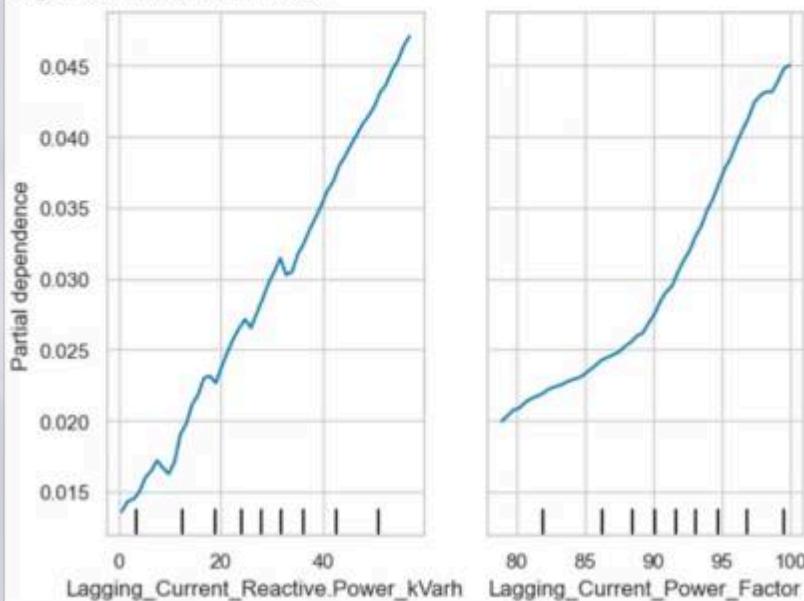
```
#from sklearn.inspection import PartialDependenceDisplay
import matplotlib.pyplot as plt

features = ['Lagging_Current_Reactive.Power_kVarh', 'Lagging_Current_Power_Factor']

# Set the figure size before creating the display
plt.figure(figsize=(10, 4))

# Create and plot the partial dependence
display = PartialDependenceDisplay.from_estimator(
    rf_model, X_test, features, kind='average', grid_resolution=50,
    feature_names=X_co2.columns
)

plt.tight_layout()
plt.show()
```



The PDPs reveal system inflection zone or points where emissions begin reacting more sharply to internal dynamics. The shape of these lines tells a story of thresholds, control transitions, and the subtle nonlinearities hiding in everyday operations.

### Lagging\_Current\_Reactive.Power\_kVarh

- Shows a stepwise increase up to ~35, followed by a sharp upward inflection.
- Suggests there's a threshold or tipping point in reactive power beyond which CO<sub>2</sub> emissions escalate.
- May reflect systems shifting into less efficient states or triggering compensatory responses after certain load levels.

### Lagging\_Current\_Power\_Factor

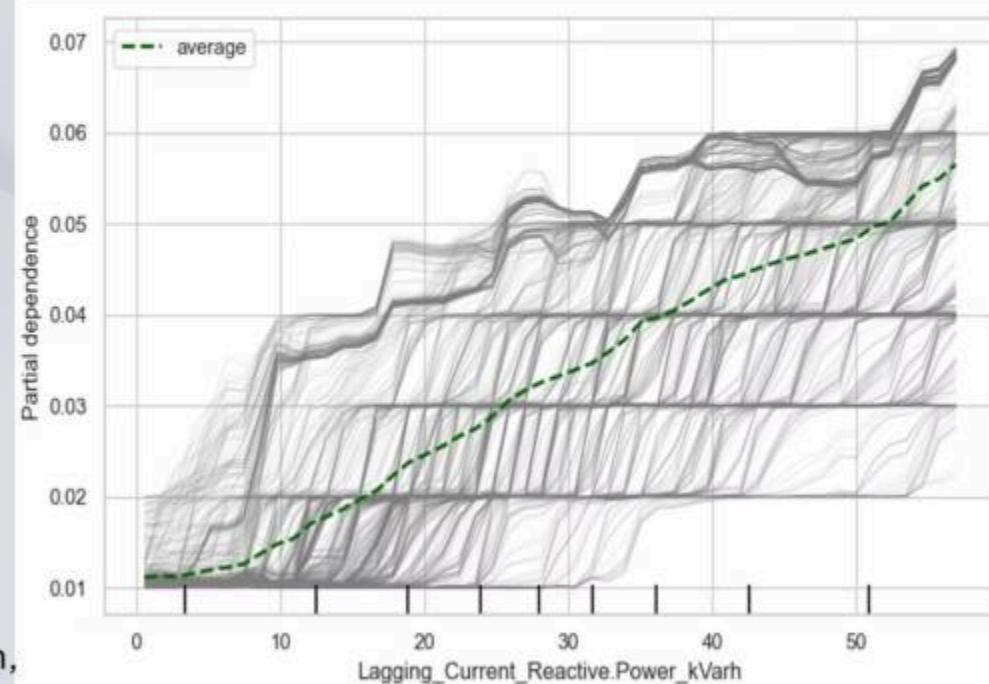
- Begins with a smooth, curved incline, possibly reflecting improving efficiency.
- Around 90–97.5, the slope steepens and then flattens momentarily before resuming upward toward 100.
  - Indicates nonlinear behavior where efficiency gains in this range correlate with increased CO<sub>2</sub>.
  - This is potentially due to load balancing or control system effects.

## ICE PLOT OF LAGGING Current Reactive Power KVarh

```
#What ICE plots show -While the PDP Line shows the average effect of a feature (like a smoothed curve),  
#ICE plots reveal how different samples react to changing just that one feature.  
#Like multiple personal stories stacked behind the summary.  
#trying it for my top feature in the model -lagging_Current_Reactive.Power_kVarh  
from sklearn.inspection import PartialDependenceDisplay  
import matplotlib.pyplot as plt  
  
# Create figure with desired size first  
fig, ax = plt.subplots(figsize=(8, 5))  
  
# Individual curves for each sample + PDP overlay  
PartialDependenceDisplay.from_estimator(  
    rf_model, X_test,  
    features=['lagging_Current_Reactive.Power_kVarh'],  
    kind='both', # both PDP and ICE  
    ice_lines_kw={'alpha': 0.2, 'color': 'gray'}, # subtle ICE lines  
    pd_line_kw={'color': 'darkgreen', 'linewidth': 2}, # clear PDP curve  
    feature_names=X_co2.columns, # Make sure this variable is defined  
    grid_resolution=50,  
    ax=ax) # Pass the axes object instead of figsize  
)  
plt.tight_layout()  
plt.show()
```

- The green dashed line, though not perfectly straight, shows a steady upward trend and indicating that as reactive power rises, predicted CO<sub>2</sub> emissions tend to increase.
- That non-linearity, gentle at first, then sharper also suggests there may be a threshold or tipping point in system behavior, where inefficiencies start to manifest more clearly.
- The strands themselves highlight sample-level variation, reflecting how the model's response isn't uniform but influenced by interaction with other features (like power factor or environmental conditions).

*The strands look almost like threads being pulled upward with each one a system snapshot, rising with the tide of reactive power.*



```

#Plot ICE by group
# Loop through humidity or season groups and generate ICE plots for each group separately for each.
#fun stuff

from sklearn.inspection import PartialDependenceDisplay
import matplotlib.pyplot as plt
import pandas as pd

# First, let's ensure we're using the exact same feature set that the model was trained on.
# Assuming rf_model was trained on X_col
feature_names = rf_model.feature_names_in_ # Get the exact feature names used during training

for humidity_bin in df_nonzero['humidity_bin'].unique(): # Changed variable name from 'season' to 'humidity'
    subset = df_nonzero[df_nonzero['humidity_bin'] == humidity_bin] # Using humidity_bin here

    # Create X_subset with exactly the same features in the same order as used during training
    X_subset = pd.DataFrame(subset[feature_names].values, columns=feature_names)
    y_subset = subset['CO2_restored'] # or y_CO2 if you're comparing

    print(f"\nICE plot for humidity: {humidity_bin}") # Now humidity_bin is defined

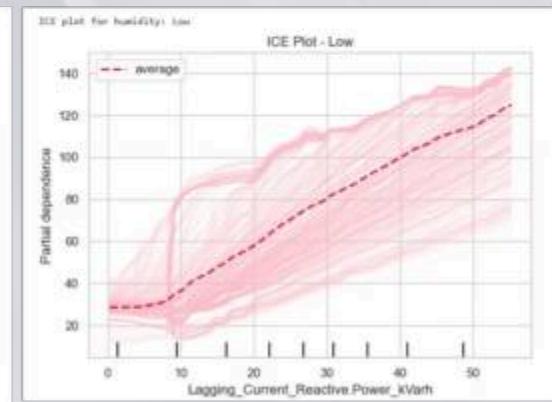
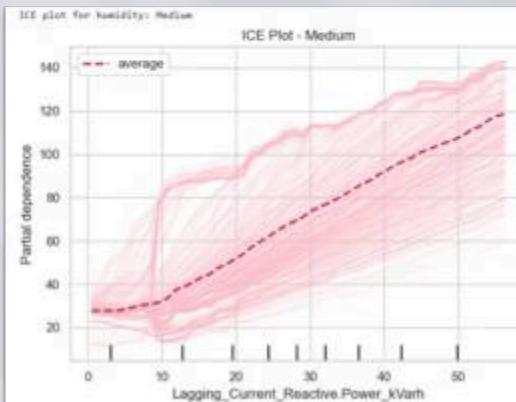
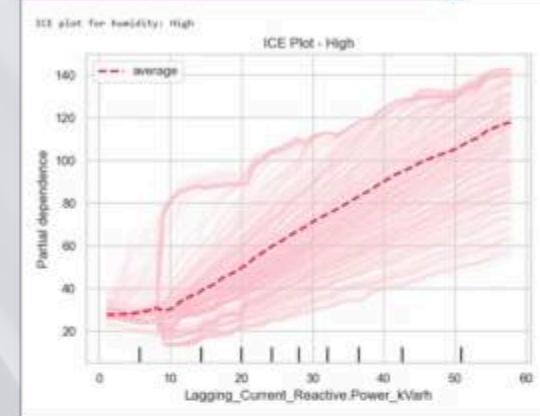
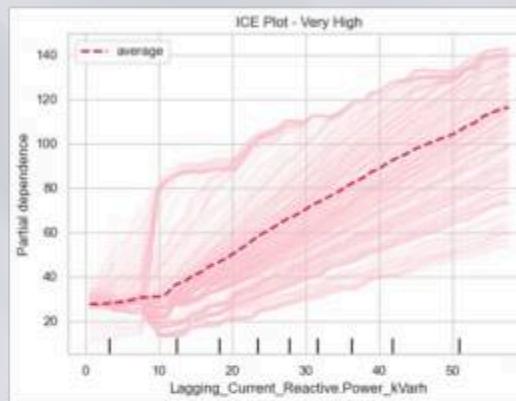
    # Use feature index instead of name to avoid any naming issues
    feature_idx = list(feature_names).index('Lagging_Current_Reactive_Power_kVarh')

    PartialDependenceDisplay.from_estimator(
        rf_model,
        X_subset,
        feature_idx, # Use index instead of name
        kind='both',
        ice_lines_kw={'alpha': 0.2, 'color': 'pink'},
        pd_line_kw={'color': 'crimson', 'linewidth': 2},
        grid_resolution=50
    )
    plt.title(f"ICE Plot - ({humidity_bin})") # Now humidity_bin is defined

```

- Low humidity = stable system response** → possibly easier control of inductive loads or better ventilation dynamics.
- Medium humidity = transitional behavior** → system may start exhibiting variability or environmental sensitivity.
- High/Very High = overlapping curves** → increased moisture could dampen reactive power's distinct influence, creating **response flattening or saturation**.

## ICE PLOTS of CO2 Looping through Humidity Bins with Feature Lagging Current Reactive Power KVarh



*Analogy: It is like tuning an instrument on under low humidity, the strings respond cleanly. However, as moisture builds, the notes blur and overlap.*

## RANDOM FOREST (Usage\_kWh without co<sub>2</sub>)

### COMPARISON: Adding and Reducing Variables

```
# Rainforest Model Usage_kWh without CO2 and added Leading_Current_Reactive_Power_kVarh

# Model for kWh with additional variable
X_kwh = df_nonzero[
    'temperature_2m (°C)',
    'relative_humidity_2m (%)',
    'Lagging_Current_Reactive_Power_kVarh',
    'Leading_Current_Power_Factor',
    'Leading_Current_Power_Factor',
    'Leading_Current_Reactive_Power_kVarh']
y_kwh = df_nonzero['Usage_kWh']

# Running same model setup as before , so beware of potential error of missing something

# Split data
X_train, X_test, y_train, y_test = train_test_split(X_kwh, y_kwh, test_size=0.2, random_state=42)

# Create and train model
rf_model = RandomForestRegressor(n_estimators=100, max_depth=None, random_state=42)
rf_model.fit(X_train, y_train)

# Run predictions and evaluate
y_pred = rf_model.predict(X_test)
print('R2 Score:', r2_score(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

# Feature importance
importance_df = pd.DataFrame({
    'feature': X_kwh.columns, # Using X_kwh columns
    'importance': rf_model.feature_importances_
})
print(importance_df.sort_values('importance', ascending=False))

R2 Score: 0.9983005599560109
RMSE: 1.0965889615560434
```

#### Feature Importance:

	feature	importance
2	Lagging_Current_Reactive_Power_kVarh	0.662485
3	Lagging_Current_Power_Factor	0.318633
4	Leading_Current_Power_Factor	0.014047
0	temperature_2m (°C)	0.002462
5	Leading_Current_Reactive_Power_kVarh	0.001795
1	relative_humidity_2m (%)	0.000578

```
# Streamlined version of Rainforest - only keeping stronger contributors

# kwh version streamlined
X_kwh = df_nonzero[
    'Lagging_Current_Reactive_Power_kVarh',
    'Lagging_Current_Power_Factor',
    'Leading_Current_Power_Factor'
] # Removed temperature, humidity, and Leading_Current_Reactive_Power_kVarh
y_kwh = df_nonzero['Usage_kWh']

# Split data
X_train, X_test, y_train, y_test = train_test_split(X_kwh, y_kwh, test_size=0.2, random_state=42)

# Create and train model
rf_model = RandomForestRegressor(n_estimators=100, max_depth=None, random_state=42)
rf_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred = rf_model.predict(X_test)
print('R2 Score:', r2_score(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

# Feature importance
importance_df = pd.DataFrame({
    'feature': X_kwh.columns, # Using X_kwh columns
    'importance': rf_model.feature_importances_
})
print('\nFeature Importance')
print(importance_df.sort_values('importance', ascending=False))

# Cross-validation
scores = cross_val_score(rf_model, X_kwh, y_kwh, cv=5)
print('\nCross-validation scores:', scores)
print('Mean CV score:', scores.mean())

R2 Score: 0.969565780686782
RMSE: 0.17464776994087217
```

#### Feature Importance:

	feature	importance
0	Lagging_Current_Reactive_Power_kVarh	0.639647
1	Lagging_Current_Power_Factor	0.291156
2	Leading_Current_Power_Factor	0.069197

Cross-validation scores: [0.94688156 0.95603223 0.96954724 0.97113622 0.96314634]  
 Mean CV score: 0.961348718241472

## RANDOM FOREST (CO<sub>2</sub> without Usage\_kWh)

### Streamlining Variables

After seeing how well the streamlined model performed for Energy Usage (kWh), I couldn't resist revisiting the CO<sub>2</sub> Random Forest model with fresh eyes. Interestingly, its structure also hinted at a 'streamlined feel' that was dominated by a few core features with consistent behavior across system states

```
# Streamlined version - only keeping stronger contributors

X_co2 = df_nonzero[['
    'Lagging_Current_Reactive.Power_kVArh',
    'Lagging_Current_Power_Factor',
    'Leading_Current_Power_Factor'
]] # Removed temperature, humidity, and Leading_Current_Reactive.Power_kVArh
y_co2 = df_nonzero['yj_CO2']

# Split data
X_train, X_test, y_train, y_test = train_test_split(X_co2, y_co2, test_size=0.2, random_state=42)

# Create and train model
rf_model = RandomForestRegressor(n_estimators=100, max_depth=None, random_state=42)
rf_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred = rf_model.predict(X_test)
print('R2 Score:', r2_score(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

# Feature importance
importance_df = pd.DataFrame({
    'feature': X_co2.columns,
    'importance': rf_model.feature_importances_
})
print('\nFeature Importance:')
print(importance_df.sort_values('importance', ascending=False))

# Cross-validation
scores = cross_val_score(rf_model, X_co2, y_co2, cv=5)
print('\nCross-validation scores:', scores)
print('Mean CV score:', scores.mean())
```

```
R2 Score: 0.969565780686782
RMSE: 0.17464776994087217

Feature Importance:
              feature  importance
0  Lagging_Current_Reactive.Power_kVArh  0.639647
1          Lagging_Current_Power_Factor  0.291156
2          Leading_Current_Power_Factor  0.069197

Cross-validation scores: [0.94688156 0.95603223 0.96954724 0.97113622 0.96314634]
Mean CV score: 0.961348718241472
```

Comparison and Overview on Next Slide



## RANDOM FOREST (Usage\_kWh without CO<sub>2</sub>) COMPARISON: Adding and Reducing Variables

# Exploring Internal System Behavior

This streamlined model shows how reactive power and system efficiency are the key engines of kWh usage. Environmental variables had minimal influence, highlighting the system's internally driven behavior. Including Leading Reactive Power added a dimension but even without it, the model held strong, affirming the resilience of the core predictors.

## Performance & Feature Impact

Version	R <sup>2</sup> Score	RMSE	Top Features
<b>Full Model (6 features)</b>	0.9983	1.0966	Lagging Reactive Power, Lagging Power Factor
<b>Streamlined (3 features)</b>	0.9978	1.2476	Same top 3, stronger individual contributions

- **Leading Reactive Power** had a small but non-zero impact in the full model, reinforcing its subtle influence.
- Streamlining preserved predictive power with **minimal loss**, proving the value of focusing on operational features.
- **Cross-validation scores** confirm robustness and generalizability across samples.

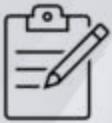
## RANDOM FOREST (CO<sub>2</sub> without Usage\_kWh) Streamlining Variables

### Comparative Summary

	Full Model (6 features)	Streamlined Model (3 features)
<b>R<sup>2</sup> Score</b>	0.9804	0.9696
<b>RMSE</b>	0.00176	0.17465
<b>Top Feature</b>	Lagging Reactive Power (65.8%)	Lagging Reactive Power (63.96%)
<b>Secondary Features</b>	Lagging PF (31.2%), Leading PF (1.6%)	Lagging PF (29.12%), Leading PF (6.92%)
<b>Environmental Variables</b>	Present (humidity, temp, etc.) – minimal impact	Excluded – focus purely on core electrical behavior
<b>Added Variable</b>	Leading Reactive Power (small but above humidity)	Not included
<b>Cross-Validation Mean</b>	Not specified	0.9613

- The full model offers superior accuracy and lower error, but confirms that most prediction strength still comes from internal system metrics.
- The streamlined model holds its own, showing minimal loss of performance while enhancing simplicity and interpretability.
- Adding environmental variables and leading reactive power helped reveal their modest but non-negligible roles, especially for deeper insights and SHAP plots.
- Both models affirm that Lagging Reactive Power and Power Factor metrics are the powerhouse predictors.

## Forecasts for VECM - Long-Term Balance, Short-Term Pulse



With time series and seasonal structure showing strong influence across models, I revisited the VECM forecast to observe how emissions and energy usage unfold over time. The following plots reflect not only directional alignment but also systemic behavior across hours and operational states, underscoring the importance of embedding temporal intelligence into future models.

VECM Model Summary (k_ar_diff=2):						
Det. terms outside the count. relation & lagged endog. parameters for equation yj_CO2						
coef	std err	z	P> z	[0.025	0.975]	
const	-9.4366	0.120	-78.734	0.000	-9.672	-9.282
exog1	-0.0070	0.000	-19.889	0.000	-0.008	-0.006
exog2	-0.0002	0.000	-0.849	0.396	-0.001	0.000
exog3	0.0512	0.000	114.001	0.000	0.050	0.052
exog4	0.0762	0.001	85.400	0.000	0.074	0.078
exog5	0.0251	0.001	38.813	0.000	0.024	0.027
exog6	0.3986	0.011	36.526	0.000	0.377	0.420
L1.yj_CO2	-0.4560	0.011	-42.324	0.000	-0.477	-0.435
L1.Usage_kWh	0.0142	0.000	32.939	0.000	0.013	0.015
L2.yj_CO2	-0.2509	0.011	-23.273	0.000	-0.272	-0.238
L2.Usage_kWh	0.0073	0.000	17.336	0.000	0.007	0.008
Det. terms outside the count. relation & lagged endog. parameters for equation Usage_kWh						
coef	std err	z	P> z	[0.025	0.975]	
const	-291.2151	2.294	-126.972	0.000	-295.710	-288.728
exog1	-0.1888	0.007	-28.016	0.000	-0.201	-0.175
exog2	-0.0866	0.004	-1.782	0.075	-0.014	0.001
exog3	1.5995	0.009	186.199	0.000	1.583	1.616
exog4	2.4414	0.017	142.960	0.000	2.408	2.475
exog5	0.6742	0.016	43.206	0.000	0.644	0.705
exog6	4.8441	0.208	23.259	0.000	4.436	5.252
L1.yj_CO2	3.5184	0.206	17.064	0.000	3.114	3.923
L1.Usage_kWh	-0.2236	0.008	-27.099	0.000	-0.240	-0.207
L2.yj_CO2	1.2505	0.206	6.063	0.000	0.846	1.655
L2.Usage_kWh	-0.0957	0.008	-11.008	0.000	-0.112	-0.080
Loading coefficients (alpha) for equation yj_CO2						
coef	std err	z	P> z	[0.025	0.975]	
ec1	-0.1732	0.001	-147.500	0.000	-0.176	-0.171
Loading coefficients (alpha) for equation Usage_kWh						
coef	std err	z	P> z	[0.025	0.975]	
ec1	-4.7168	0.022	-209.859	0.000	-4.761	-4.673
Cointegration relations for loading-coefficients-column 1						
coef	std err	z	P> z	[0.025	0.975]	
beta.1	1.0000	0	0	0.000	1.000	1.000
beta.2	0.1401	0.001	168.592	0.000	0.138	0.142

### Overview of Graph

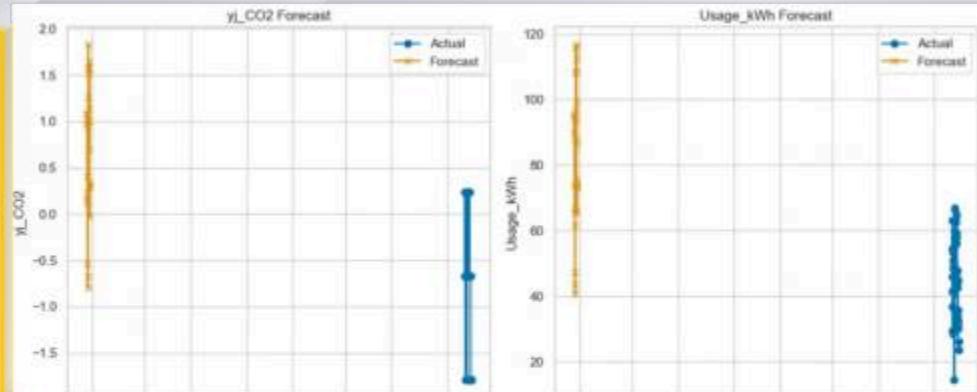
#### •Directionally accurate:

Both usage and CO<sub>2</sub> appear to move in sync between forecast and actuals.

#### •Magnitude mismatch:

The actuals dip lower and fluctuate more possibly due to seasonal or load-based factors the model didn't capture.

•**VECM nuance:** Because VECM handles long-term equilibrium and short-term dynamics, smaller mismatches may be expected unless lag order and cointegration vectors are optimized for seasonal shifts.



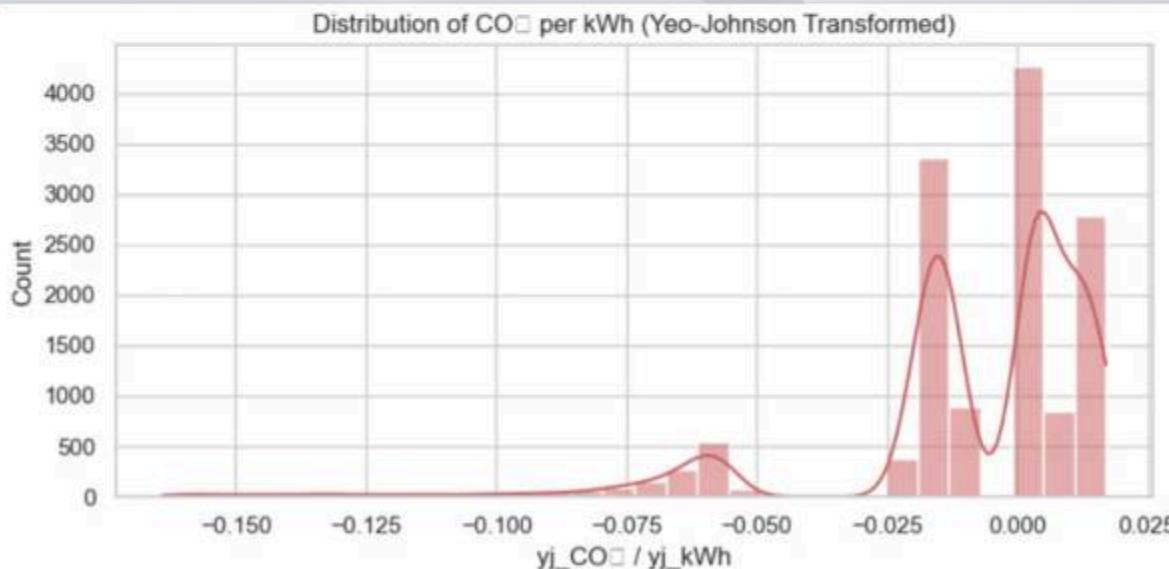
Target	Test MSE	Interpretation
yj_CO2	64.46	Moderate fit; residual spread suggests complexity or transformation effects
yj_kWh	2792.16	Larger errors; possibly due to variability in usage behavior or nonlinear load effects

These results suggest CO<sub>2</sub> and energy usage are tightly coupled in the long term but respond differently in the short term, possibly reflecting **efficiency transitions, reactive compensation, or load rank dynamics**. The strong loading on Usage\_kWh confirms its role in anchoring the system, while CO<sub>2</sub> shows a subtler drift and correction behavior.

## CO<sub>2</sub> per kWh: Yeo-Johnson Transformed Distribution

*Creation of a target variable for CO<sub>2</sub> per unit of energy (CO<sub>2</sub> / kWh) to possibly use with other modeling*

```
#Creating a CO2 per kWh target to use for modeling.  
#use transformed variables so everything is normalized  
df_nonzero['CO2_per_kWh'] = df_nonzero['yj_CO2'] / df_nonzero['Usage_kWh']  
#Exploring the distribution. Want to assess the outliers and skew and other stuff that might be going on  
import seaborn as sns  
plt.figure(figsize=(8, 4))  
sns.histplot(df_nonzero['CO2_per_kWh'], bins=30, kde=True, color='indianred')  
plt.title('Distribution of CO2 per kWh (Yeo-Johnson Transformed)')  
plt.xlabel('yj_CO2 / yj_kWh')  
plt.tight_layout()  
plt.show()
```



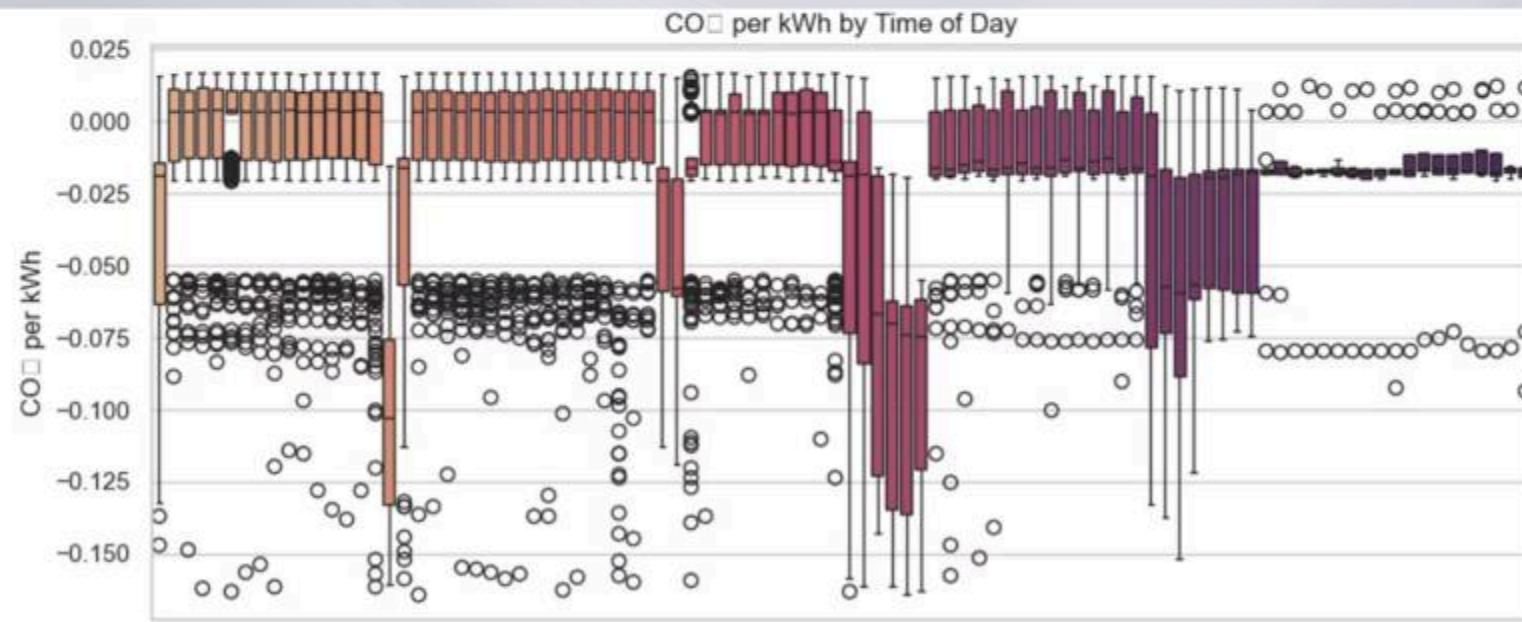
**The transformed distribution revealed several distinct behaviors:**

- A small bell-shaped cluster from ~−0.075 to +0.05 possibly representing low-load or highly efficient states.
- A tighter triad of bars around −0.025, where the middle bar peaks and could reflect system conditions where emissions efficiency stabilizes before shifting again.
- At zero, the distribution breaks symmetry: the tallest bar appears here, followed by a dip, then a partial rise . Again, suggesting a high-frequency operating zone with surrounding variability.

## Box Plot of CO<sub>2</sub>\_per\_kWh by Time of Day



This boxplot reveals how emissions intensity (CO<sub>2</sub>/kWh) fluctuates across the day. While visually dense, it exposes clear structural changes—including high variability in the early morning hours, localized dips around 9:30pm and 12:20am, and tightly grouped distributions during certain stable periods. The presence of frequent outliers suggests transient events or operational shifts worth exploring time series further.

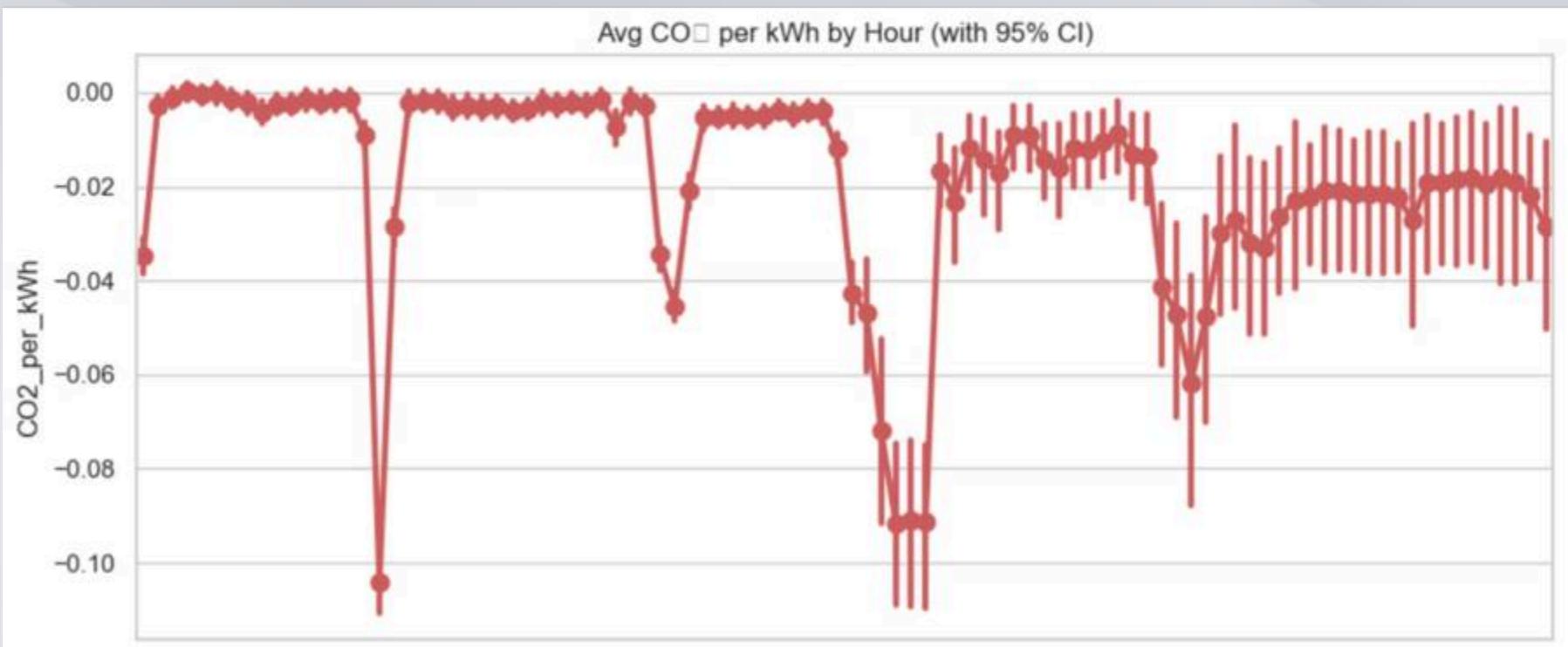


### Observations from Plot

- The **early spikes and deep whiskers** suggest emissions intensity is most unstable during late-night hours and possibly tied to reactive load conditions or transitional states.
- The **clustered shifts around 9:30pm and 12:20am** could point to control handoffs, efficiency plateaus, or environmental response changes.
- The **tiny plots post-5am** might be telling a tale and suggesting minimal variation due to possibly idle phase.

## Hourly Trend: Avg CO<sub>2</sub> per kWh with Confidence Bands

To validate the time-dependent behavior observed in the boxplots, I generated a point plot of average CO<sub>2</sub> per kWh across hours, including 95% confidence intervals. Despite some limitations in temporal granularity, the overall shape mirrors earlier findings: **higher variability in early morning hours**, subtle **dips around late night**, and pockets of stability that reflect consistent system states



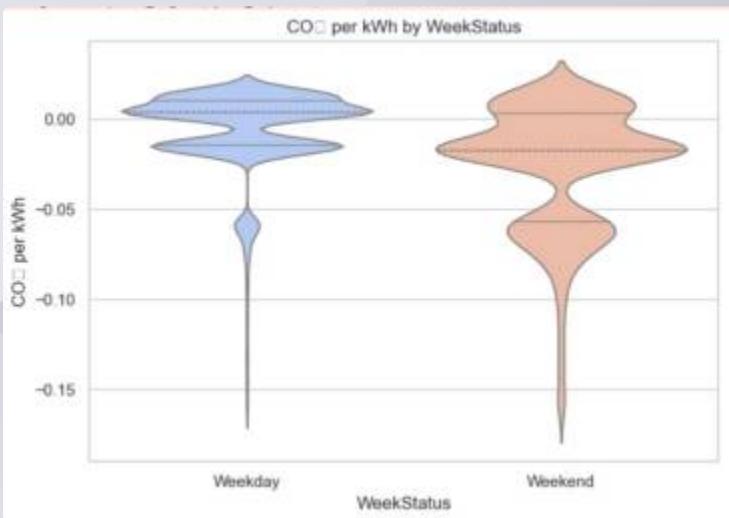


# Violin Plot of CO<sub>2</sub> per kWh Distribution by Day Type

## Exploring Weekday vs. Weekend Emission Profiles

Given the temporal sensitivity observed in both the CO<sub>2</sub> and Energy Usage (kWh) models, I expanded the analysis to uncover additional time-driven behaviors. Beyond seasonal trends and hourly fluctuations, the data hinted at possible operational patterns, such as shift transitions or downtime periods, that could impact emissions and usage metrics.

**The violin plot reveals distinct structural differences in emissions between weekdays and weekends**

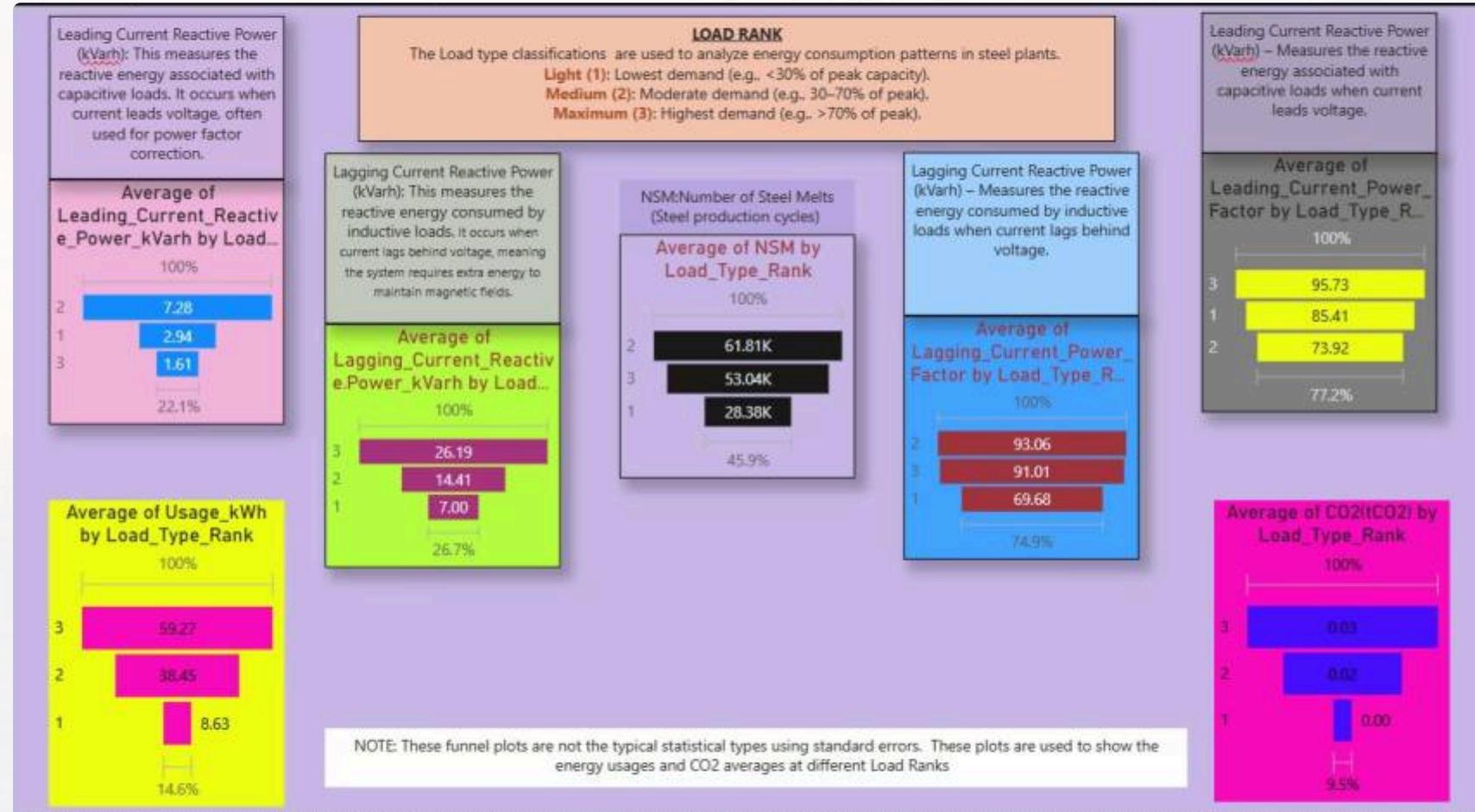


**Weekdays** show a wider top density near zero, narrowing quickly with a brief blip around -0.05, then tapering into a thin tail toward -0.15. This suggests a central concentration of operating states, with tightly regulated emissions and fewer extreme outliers.

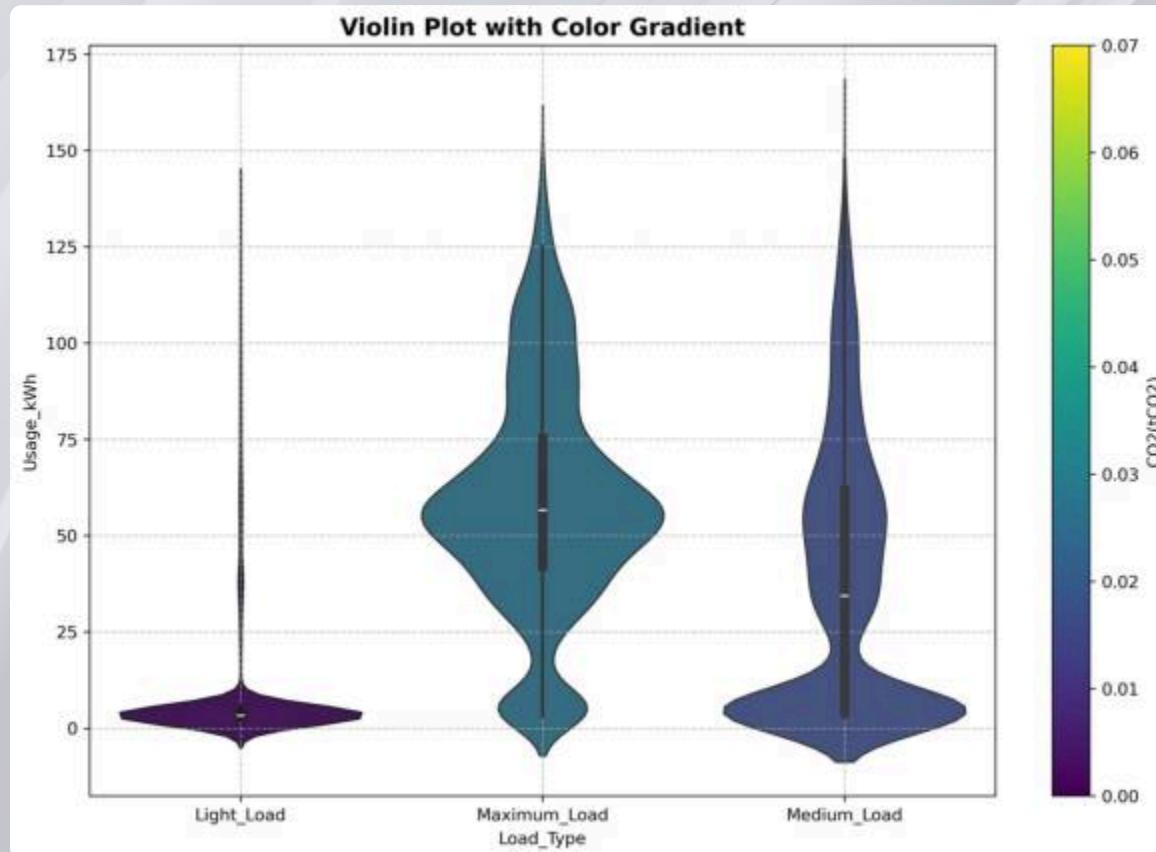
**Weekends**, in contrast, feature a slightly flatter top, followed by a broader second segment, especially near -0.05, where density thickens. Rather than tightening into a fine tail like weekdays, the weekend distribution remains more dispersed and rounded, implying diverse operational conditions or less predictable system usage.

# Overview of Load Ranks

Load ranks are reintroduced in the Random Forest models to assess their predictive impact and clarify nonlinear energy behavior.



# Violin Plot of Categorical Load vs. Energy Usage kWh



Before diving back into Random Forest with categorical load, I decided to observe how CO<sub>2</sub> emissions looks like with energy usage across load ranks.

Light Load presents a tight vertical band and modest spread, concentrated operation. Medium Load extends upward with a slightly looser curve, signaling more variability but still restrained. Max Load tells a different story: visible shifts, broader density changes, and dynamic expansion patterns suggest system transitions, constraint thresholds, or regime shifts.

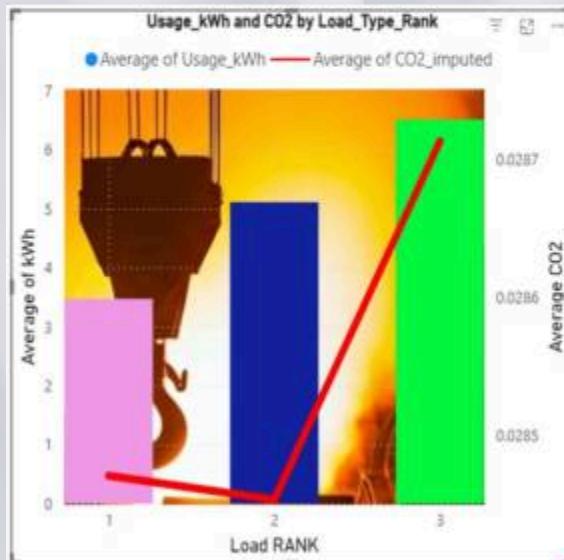
These distributional cues preview what the model will later confirm: operational load rank isn't just a label; it's a behavioral fingerprint.

## Random Forest Model with Categorical Load

### Load Rank Revisited: Operational Demand Meets Emission Behavior



Earlier in my analysis, Load Rank emerged as a key operational marker with distinguishing low, medium, and peak production states. I examined its relationship to energy consumption and emissions behavior, uncovering nonlinear shifts and efficiency transitions. Now, using Random Forest models tailored to both kWh and CO<sub>2</sub>, I return to Load Rank not just as a descriptive feature, but as a **predictive anchor** with letting the model map how systemic load translates into energy usage and environmental impact.

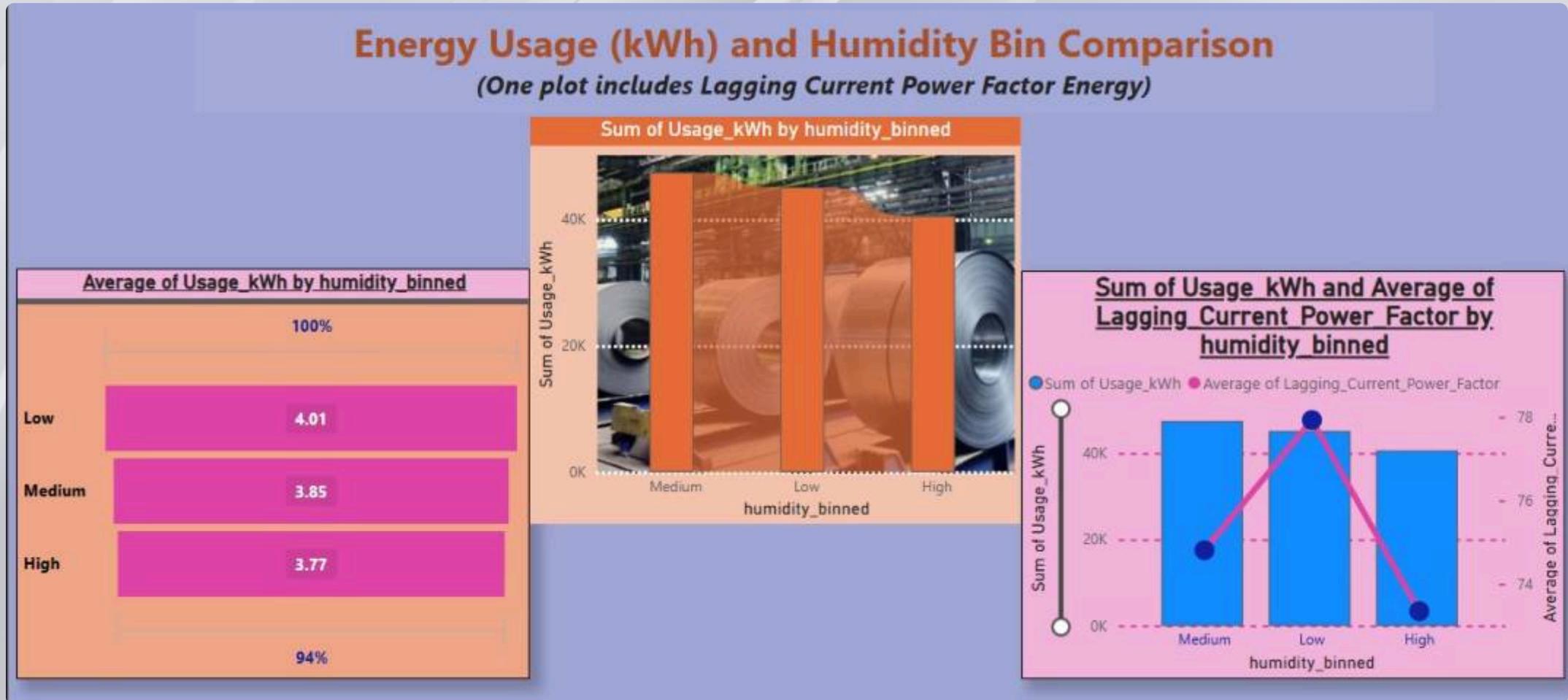


CO2 R2 Score: 0.9708050923470588	R <sup>2</sup> -CO <sub>2</sub> : ~0.971																					
CO2 RMSE: 0.17105489808566518	R <sup>2</sup> -kWh: ~0.998																					
CO2 Feature Importance:																						
<table><thead><tr><th></th><th>feature</th><th>importance</th></tr></thead><tbody><tr><td>0</td><td>Lagging_Current_Reactive_Power_kVarh</td><td>0.633427</td></tr><tr><td>1</td><td>Lagging_Current_Power_Factor</td><td>0.284339</td></tr><tr><td>2</td><td>Leading_Current_Power_Factor</td><td>0.066426</td></tr><tr><td>3</td><td>temperature_2m (°C)</td><td>0.015792</td></tr><tr><td>4</td><td>Load_Type_Medium_Load</td><td>0.000011</td></tr><tr><td>5</td><td>Load_Type_Maximum_Load</td><td>0.000005</td></tr></tbody></table>			feature	importance	0	Lagging_Current_Reactive_Power_kVarh	0.633427	1	Lagging_Current_Power_Factor	0.284339	2	Leading_Current_Power_Factor	0.066426	3	temperature_2m (°C)	0.015792	4	Load_Type_Medium_Load	0.000011	5	Load_Type_Maximum_Load	0.000005
	feature	importance																				
0	Lagging_Current_Reactive_Power_kVarh	0.633427																				
1	Lagging_Current_Power_Factor	0.284339																				
2	Leading_Current_Power_Factor	0.066426																				
3	temperature_2m (°C)	0.015792																				
4	Load_Type_Medium_Load	0.000011																				
5	Load_Type_Maximum_Load	0.000005																				
CO2 Cross-validation scores: [0.93856306 0.93809676 0.96133902 0.96507024 0.93818994]																						
CO2 Mean CV score: 0.9482518047625323																						
kWh R2 Score: 0.997779363333045																						
kWh RMSE: 1.25351547410715																						
kWh Feature Importance:																						
<table><thead><tr><th></th><th>feature</th><th>importance</th></tr></thead><tbody><tr><td>0</td><td>Lagging_Current_Reactive_Power_kVarh</td><td>0.663045</td></tr><tr><td>1</td><td>Lagging_Current_Power_Factor</td><td>0.319302</td></tr><tr><td>2</td><td>Leading_Current_Power_Factor</td><td>0.011740</td></tr><tr><td>3</td><td>temperature_2m (°C)</td><td>0.002908</td></tr><tr><td>4</td><td>Load_Type_Medium_Load</td><td>0.001680</td></tr><tr><td>5</td><td>Load_Type_Maximum_Load</td><td>0.001324</td></tr></tbody></table>			feature	importance	0	Lagging_Current_Reactive_Power_kVarh	0.663045	1	Lagging_Current_Power_Factor	0.319302	2	Leading_Current_Power_Factor	0.011740	3	temperature_2m (°C)	0.002908	4	Load_Type_Medium_Load	0.001680	5	Load_Type_Maximum_Load	0.001324
	feature	importance																				
0	Lagging_Current_Reactive_Power_kVarh	0.663045																				
1	Lagging_Current_Power_Factor	0.319302																				
2	Leading_Current_Power_Factor	0.011740																				
3	temperature_2m (°C)	0.002908																				
4	Load_Type_Medium_Load	0.001680																				
5	Load_Type_Maximum_Load	0.001324																				
kWh Cross-validation scores: [0.99297465 0.99092111 0.99599456 0.99737682 0.99540156]																						
kWh Mean CV score: 0.9945337377138947																						
note: Part of Coding was to assess load category's importance. If significant (>0.05), include in Model																						

#### Interpretive Highlights

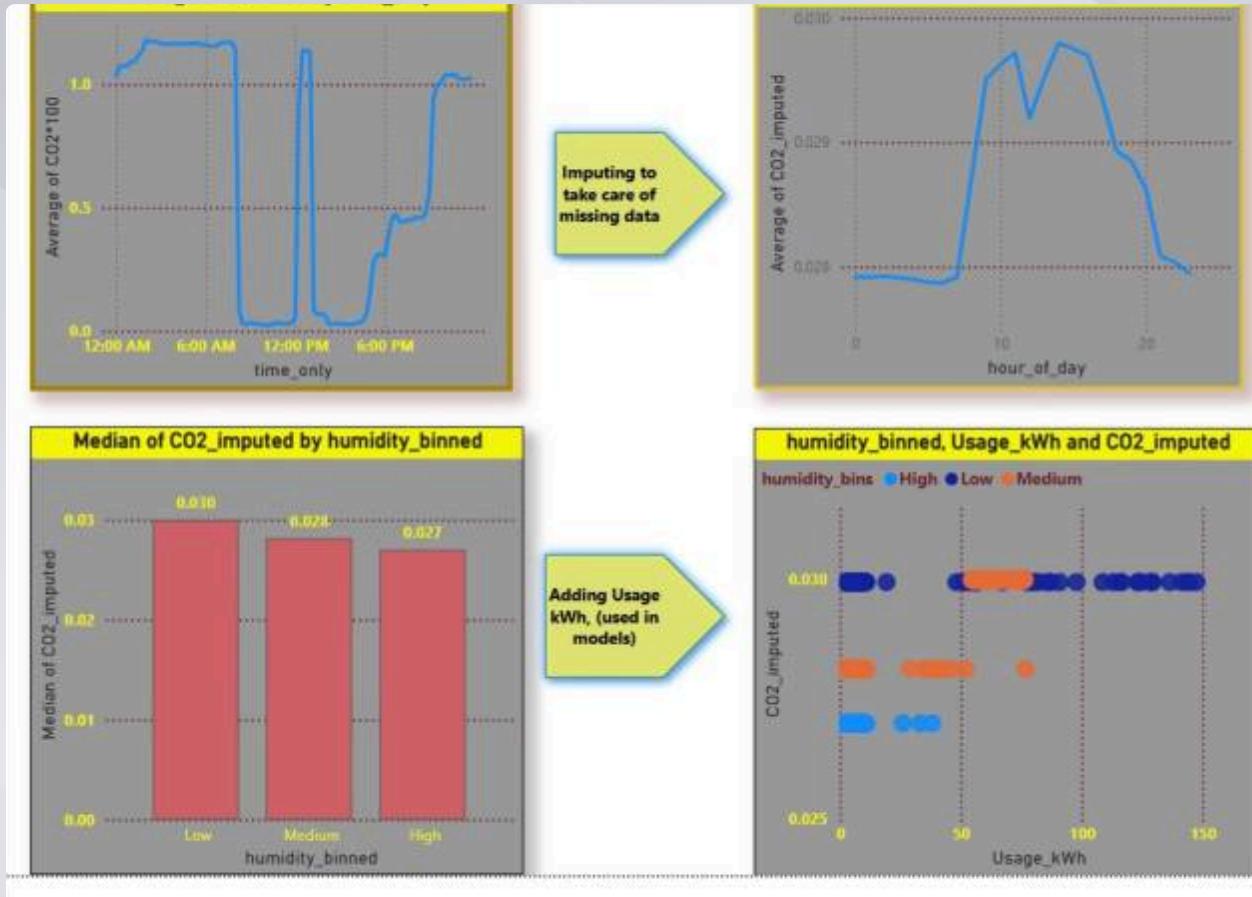
- **Reactive metrics dominate** in both models; the system's internal load dynamics are the clear drivers.
- **Temperature matters modestly**, more so for CO<sub>2</sub> than for energy and suggesting environmental influence on emissions efficiency.
- **Load Rank's low importance confirms earlier hypotheses**: while operational states influence performance qualitatively, they don't significantly shift predictive structure in this tree-based framework.
- **Cross-validation scores show excellent generalizability** with minimal overfitting and strong consistency across folds.

While categorical load patterns revealed operational structure through the RF model, system behavior doesn't unfold in isolation. Environmental conditions, particularly humidity, introduce subtle but measurable shifts in energy usage and emissions. Zooming into these climate driven regimes, it was uncovered how moisture levels reshape operational rhythm and predictive clarity.



# CO<sub>2</sub> Evaluation Continued: Time, Energy, and Humidity Interplay

## Temporal and Humidity-Driven Emission Patterns



## OVERVIEW

Actual CO<sub>2</sub> emissions reveal sharp overnight peaks and midday plateaus, while imputed values smooth these transitions. Humidity bins show clear stratification: low humidity drives high CO<sub>2</sub> across all energy levels, while high humidity suppresses emissions and caps energy usage. These patterns suggest that environmental conditions modulate system efficiency and load behavior.

# HUMIDITY: Diagnostically justified, not just Exploratory

## Dunn Test: CO<sub>2</sub> and Humidity Bins

With p-values so tiny the takeaway is loud and clear: each humidity bin's CO<sub>2</sub> distribution is significantly different from the others. Not just statistically significant, but celestially significant if that existed. This validates what a Kruskal-Wallis test hinted at earlier: The humidity's effect isn't subtle. It appears that the imputed CO<sub>2</sub> values cluster in distinctly different ways across Low, Medium, and High humidity. Looking at the charts and graph, it also shows a trend of higher humidity yields lower CO<sub>2</sub>.

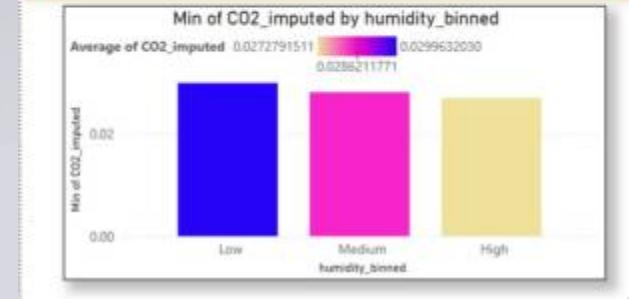
```
import scikit_posthocs as sp
import numpy as np

# Perform Dunn's test
dunn_results = sp.posthoc_dunn(df, val_col="CO2(tCO2)_imputed", group_col="humidity_binned", p_adjust="bonferroni")

print(dunn_results)
```

Dunn Test		High	Low	Medium
High	1.000000e+00	0.000000e+00	2.348857e-272	
Low	0.000000e+00	1.000000e+00	3.040094e-261	
Medium	2.348857e-272	3.040094e-261	1.000000e+00	

Computed CO <sub>2</sub> Bins by Humidity Bins CO <sub>2</sub> × 100					
humidity_binned	27-28k	28-29k	29-30k	30k+	Total
High	26251			★ 27135	27135
Low		28235	27891	29693	29693
Medium	30506		27733	33205	33205
Total	26251	30506	28235	38593	40381



# The Power of XGBRegressor Modeling....Next

To deepen my exploration of emissions and energy usage, I revisited XGBRegressor. XGBoost is a highly efficient tree-based algorithm known for its speed and predictive strength. While I've worked with XGBoost before, this pass includes a fresh build to ensure accuracy and alignment with updated transformations and system insights.

## COMPARISON OF MODELS USED FOR ANALYSIS

Model	Core Mechanism	Best For	Notes
OLS	Linear regression (ordinary least squares)	<b>Directionality &amp; baseline insights</b>	Great for quick diagnostics, identifying linear relationships
VECM	Cointegration-based system modeling	<b>Temporal equilibrium &amp; system feedback</b>	Ideal for long-run relationships and shock-response analysis
Random Forest	Ensemble tree voting (bagging)	<b>Robust prediction &amp; feature interpretability</b>	Excellent with SHAP, ICE, PDPs for uncovering nonlinear drivers
XGBRegressor	Gradient-boosted decision trees (iterative error correction)	<b>High-performance modeling</b>	Captures subtle interactions, excels in tuned prediction tasks

# XGBoost Model with log\_kWh and Temporal Variables



Initially, my exploration with XGBoost steered me away from log-transformed kWh. I was focused on minimizing complexity and leaning into raw system behavior. But as the modeling deepened, and in my quest to stabilize variance and reduce distributional skewness, I felt compelled to circle back. The log\_kWh, though once sidelined, offered the potential for behavioral clarity and improved interpretability. And honestly, I'm also just curious, which is why I want to give low humidity another look as well.

```
Python Code for Refined Model (v2d) - keeping track and saving...
From sklearn import XGBRegressor
From sklearn.metrics import r2_score, mean_squared_error
From sklearn.model_selection import train_test_split
From sklearn.preprocessing import StandardScaler
Import pandas as pd
Import numpy as np

# Load dataset
df = pd.read_csv(r'C:/Users/mash/Desktop/steel_weather_data_structured_v7.csv')
print("Dataset shape:", df.shape)
print("Low humidity rows?", len(df[df['humidity_low'] == 'Low']))

# Create log_kWh feature
df['log_kWh'] = np.log1p(df['Usage_kWh']) # (log1p to handle zeros)

# Define features (include log_kWh, exclude CO2_lag_1 to boost Usage_kWh)
features = [
    'log_kWh', 'temperature_2m (°C)', 'relative_humidity_2m (%)',
    'hour_of_day', 'is_weekend'
]
day_of_week_cols = [col for col in df.columns if col.startswith('day_of_week_')]
features.extend(day_of_week_cols)
print("Features used:", features)

# Low humidity subset
df_low = df[df['humidity_low'] == 'Low'].copy()
X_low = df_low[features]
y_low = df_low['CO2_imputed']
# Check feature statistics
print("Feature variances:", X_low.var())

# Scale features
scaler = StandardScaler()
X_low_scaled = scaler.fit_transform(X_low)
X_low_scaled = pd.DataFrame(X_low_scaled, columns=features)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X_low_scaled, y_low, test_size=0.2, random_state=42)

# Train model (fixed parameters)
model_low = XGBRegressor(
    n_estimators=2000, max_depth=1, learning_rate=0.1,
    min_child_weight=1, subsample=0.8, colsample_bytree=0.8, random_state=42
)
model_low.fit(X_train, y_train)

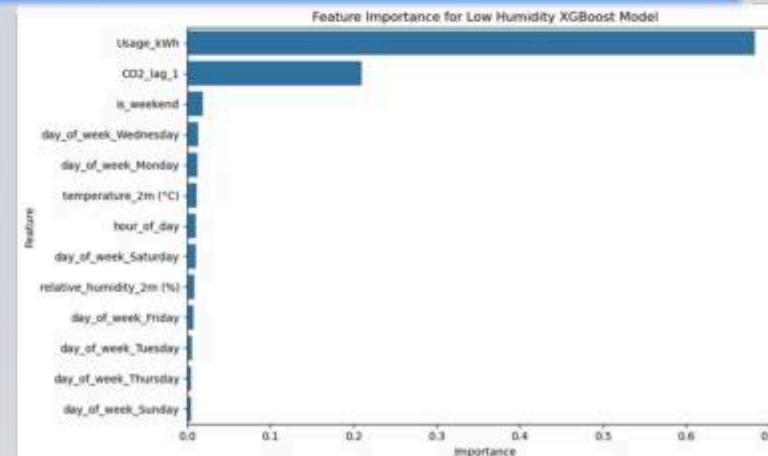
# Evaluate
y_train_pred = model_low.predict(X_train)
y_test_pred = model_low.predict(X_test)
print("Low Humidity Train R^2:", r2_score(y_train, y_train_pred))
print("Low Humidity Test R^2:", r2_score(y_test, y_test_pred))
print("Low Humidity Train MSE:", mean_squared_error(y_train, y_train_pred))
print("Low Humidity Test MSE:", mean_squared_error(y_test, y_test_pred))

# Save feature importance
importances = pd.DataFrame({'feature': features, 'Importance': model_low.feature_importances_})
print("Feature importance:")
print(importances.sort_values('Importance', ascending=False))
importances.to_csv(r'C:/Users/mash/Desktop/feature_importance_low_xgbtest_v2d.csv', index=False)
print("Feature importance saved to feature_importance_low_xgbtest_v2d.csv")
```

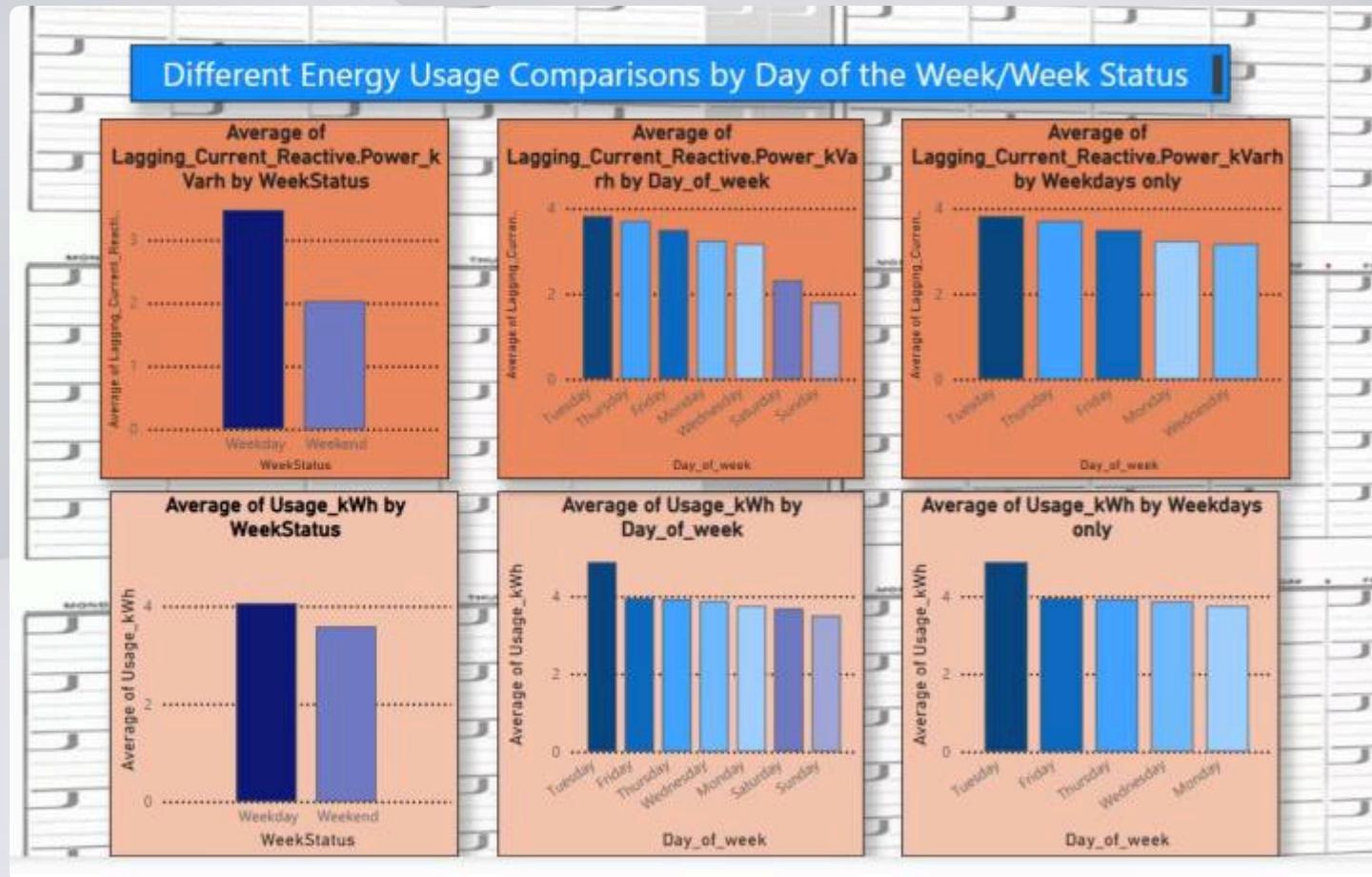
## Overview and Thoughts

- **Excellent predictive strength:**  $R^2 \sim 0.989$  on both train and test sets with microscopic MSE
- **Log transformation stabilizes variance and normalizes skew:** Given earlier diagnostics, this aligns beautifully with the system's distributional trait which is especially helpful under low humidity regimes.
- **Feature importance reveals subtle behavioral structure:** While log\_kWh dominates (~ 88.5%), variables such as, is\_weekend, hour\_of\_day, and day\_of\_week\_Sunday, still contribute meaningfully and is telling a story about temporal and operational rhythm beneath the energy curve.

Low Humidity Train  $R^2: 0.9890803090115042$   
Low Humidity Test  $R^2: 0.9897590107075699$   
Low Humidity Train MSE: 1.6371012115759553e-07  
Low Humidity Test MSE: 1.44710084974806e-07



# Temporal Transitions in Energy: A Look at Weekday vs. Weekend Trends



Not surprising, Weekdays yield more Energy usage and Lagging Current Reactive Power. Surprising, Tuesday leads in energy demand. Whether it's operational strategy or logging cadence, this day may carry the plant's productive backbone.

## This Tuesday Spike could indicate:

- **Early-week ramp-up behavior** following Monday resets or system checks.
- **Production prioritization**, where Tuesday holds heavier task loads before midweek smoothing.
- Or even **data rhythm quirks**, where full shift logs or uninterrupted measurement cycles land most clearly on Tuesdays

```

from xgboost import XGBRegressor
from sklearn.metrics import r2_score, mean_squared_error
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import pandas as pd
import numpy as np

# Load dataset
df = pd.read_csv(r'C:/Users/arsch/OneDrive/Desktop/steel_weather_data_restored_v7.csv')
print("Dataset shape:", df.shape)
print("Low Humidity count:", len(df[df['humidity_binned'] == 'Low']))

# Create log_kWh feature
df['log_kWh'] = np.log10(df['Usage_kWh'])

# Define features (include CO2_lag_1 with minimal others)
features = ['log_kWh', 'CO2_lag_1', 'is_weekend', 'day_of_week_Tuesday']
print("Features used:", features)

# Low Humidity subset
df_low = df[df['humidity_binned'] == 'Low'].copy()
X_low = df_low[features]
y_low = df_low['CO2_imputed']

# Check feature statistics
print("Feature variances:", X_low.var())

# Scale features
scaler = StandardScaler()
X_low_scaled = scaler.fit_transform(X_low)
X_low_scaled = pd.DataFrame(X_low_scaled, columns=features)
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X_low_scaled, y_low, test_size=0.2, random_state=42)

# Train model (continued with early stopping)
model_low = XGBRegressor(
    n_estimators=3000, max_depth=10, learning_rate=.001,
    min_child_weight=10, subsample=0.9, colsample_bytree=0.6, random_state=42,
    early_stopping_rounds=100, eval_metric='rmse'
)
model_low.fit(
    X_train, y_train,
    eval_set=(X_test, y_test),
    verbose=False
)

# Evaluate
y_train_pred = model_low.predict(X_train)
y_test_pred = model_low.predict(X_test)
print("Low Humidity Train R^2:", r2_score(y_train, y_train_pred))
print("Low Humidity Test R^2:", r2_score(y_test, y_test_pred))
print("Low Humidity Train MSE:", mean_squared_error(y_train, y_train_pred))
print("Low Humidity Test MSE:", mean_squared_error(y_test, y_test_pred))

# Save Feature importance
importances = pd.DataFrame({'Feature': features, 'Importance': model_low.feature_importances_})
print("Feature Importance:")
print(importances.sort_values('Importance', ascending=False))
importances.to_csv(r'C:/Users/arsch/OneDrive/Desktop/feature_importance_low_xgboost_v27.csv', index=False)
print("Feature importance saved to feature_importance_low_xgboost_v27.csv")

```

## XGBoost Forest Model with log\_kWh With limited Temporal Variables and Lag1

The previous XGBoost models (not all included) with log\_kWh under low humidity conditions has shown some great strength. I selected some specific variables, like Tuesday shown in graph, that have shown past significance from both the models and other analysis.

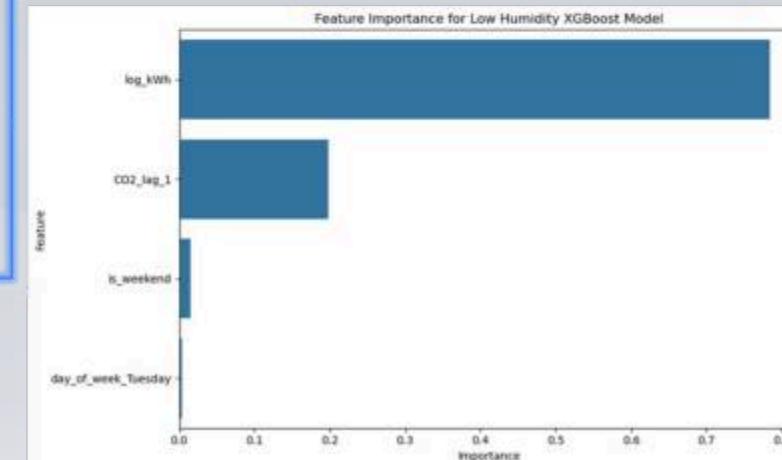
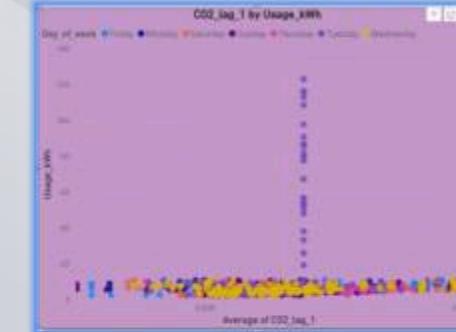
Low Humidity Train R<sup>2</sup>: 0.9865466799168923  
 Low Humidity Test R<sup>2</sup>: 0.990629537339741  
 Low Humidity Train MSE: 2.0169477900957362e-07  
 Low Humidity Test MSE: 1.3240912660866145e-07

### Overview and Thoughts

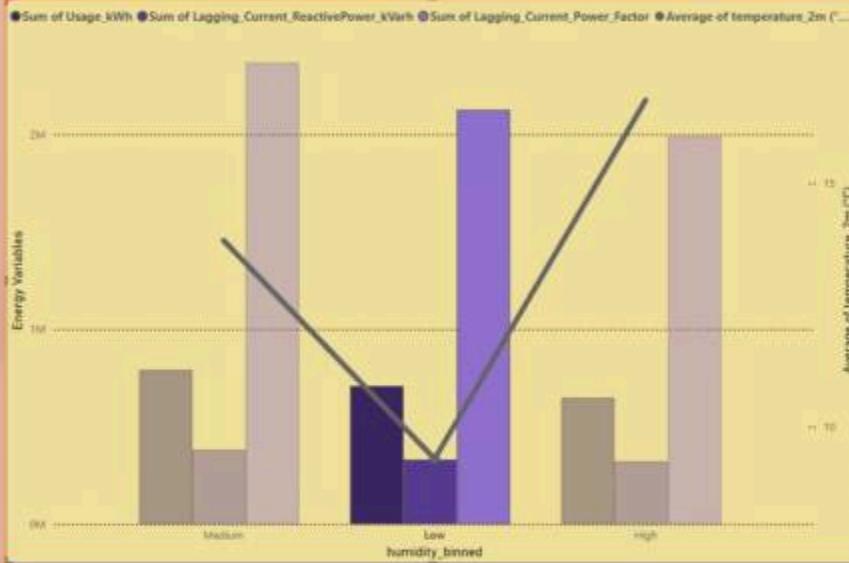
- This model hits nearly perfect R<sup>2</sup> with just four features which is proof the system's behavioral core is under dry conditions.
- The inclusion of CO2\_lag\_1 shows emission inertia matters, even when modeling energy directly.
- Temporal features like Tuesday and weekend shift model behavior subtly, revealing rhythmic patterns in operation.
- Usage doesn't just flow, it pulses..

### Feature variances:

log_kWh	1.420569
CO2_lag_1	0.000015
is_weekend	0.217252
day_of_week_Tuesday	0.097780



## Variable Benchmark: A Comparative Modeling Wrap-Up (Almost)



```
# Prepare data for CO2 prediction
features = ['Usage_kWh', 'Lagging_Current_Reactive.Power_kVarh', 'Lagging_Current_Power_Factor',
            'temperature_2m (°C)', 'relative_humidity_2m (%)']
X = df[features]
y = df['CO2(tCO2)']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Nearing the conclusion of my predictive analysis, I constructed a common feature set by leveraging variables consistently influential across prior models. This subset includes operational signals (Usage\_kWh, Lagging Reactive Power, Power Factor) and environmental indicators (temperature and relative humidity). These variables were tested under **low humidity conditions**, a regime previously shown to sharpen predictive clarity. By applying this shared structure to OLS, Random Forest, and XGBoost models, I could evaluate:

- How each model interprets the same system behaviors.
- Whether tree-based techniques outperform linear methods in feature complexity.
- Where environmental signals fade or amplify across modeling frameworks.

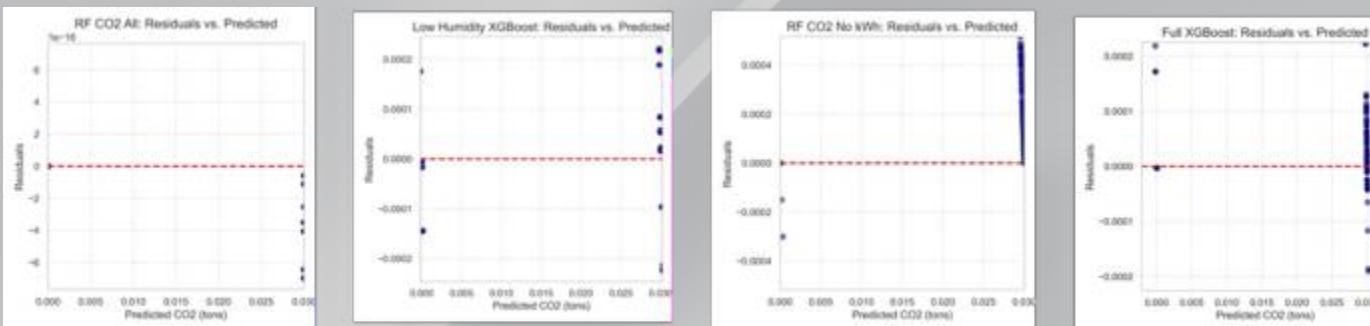
Unified Variable Suite — Model Interpretation Side-by-Side



# Model Performance Comparison : CO<sub>2</sub> Prediction

Model	R <sup>2</sup> Score	RMSE	Notes
RF (All Variables)	<b>0.9985</b>	<b>0.0004362</b>	Tight error spread, minimal outliers — robust baseline
XGBoost (Low Humidity)	<b>0.9944</b>	<b>0.0009868</b>	Best Q-Q fit, fewer outliers — strong under dry regime
Full XGBoost	<b>0.9951</b>	<b>0.0007801</b>	Slight right-tail curvature, moderate outlier density
RF (No Usage_kWh)	<b>0.9944</b>	<b>0.0008356</b>	Broader error spread, usage omission shows loss of signal

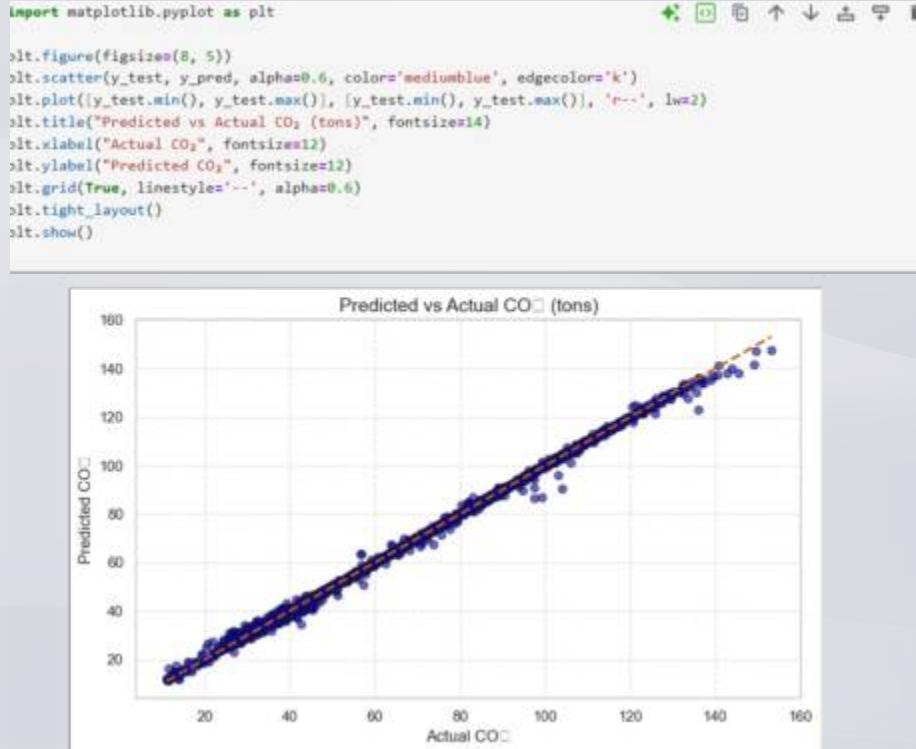
- Best Fit:** RF (All Variables) wins on raw accuracy, but **Low Humidity XGBoost** shines in interpretability and distribution alignment.
- Most Degraded Fit:** RF without Usage\_kWh continues to confirm usage is a core system signal.



**Plot 1: RF (All Variables)** shows symmetrical error with two edge outliers, indicating strong predictive balance. **Plot 2: XGBoost (Low Humidity)** exhibits the best distribution fit with minimal deviation from the normal curve. **Plot 3: RF (No Usage\_kWh)** reveals a wider error spread and structural limitation due to the omission of usage data. **Plot 4: Full XGBoost** displays right-skewed groups with possible nonlinear leakage or feature imbalance.

# The OLS Model: Final Fit

## Simplicity Earned



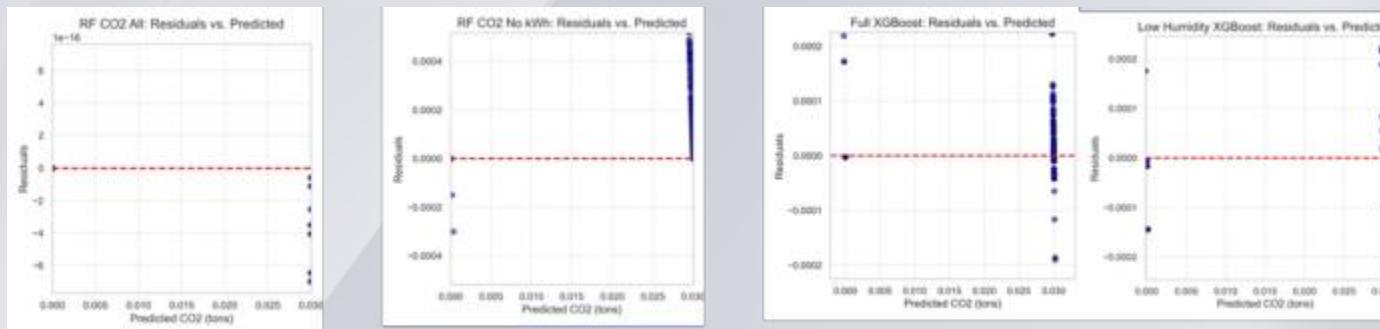
While nonlinear models captured complexity across regimes and environmental layers, the final OLS regression delivered elegant simplicity. With well-behaved inputs and a refined feature set, the predicted vs. actual CO<sub>2</sub> plot shows near-perfect alignment with each point hugging the diagonal line. This fit reflects not just statistical performance, but how far the modeling journey has come—from messy signals to interpretable clarity.

# Residual Plot Interpretation: CO<sub>2</sub> Predictions Across Models

Model	Residual Spread & Pattern	Interpretive Notes
RF (All Variables)	One point perfectly on zero; remaining outliers on the far right, extending down to ~ -0.07	Excellent predictive centering with limited spread; tight, confident
RF (No Usage_kWh)	Sparse right-side outliers; dense left-side concentration above zero	Indicates underprediction across many cases; missing dominant signal
Full XGBoost	Left: 3 dots climb upward; Right: dense cluster between -0.010 and 0.015, plus 4 key outliers	Slight overprediction tendency on right tail; shows reactive behavior
Low Humidity XGBoost	Left: 4 dots (2 near-zero); Right: more concentrated residuals above zero, few below	Most symmetrical distribution. Strong alignment under environmental control

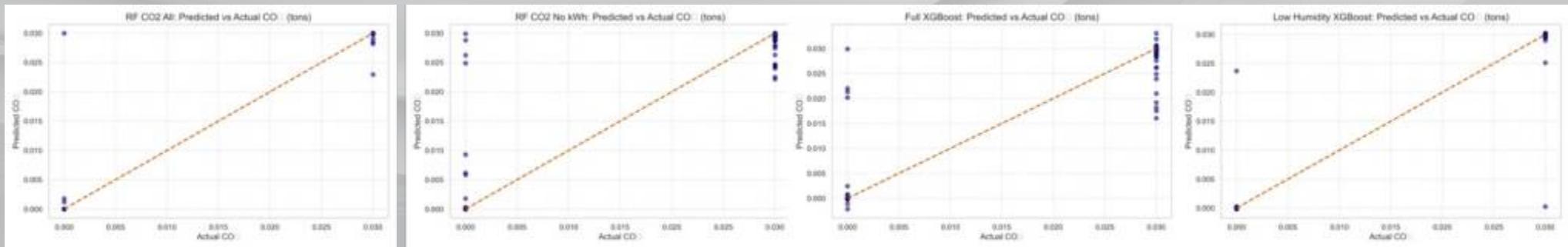
Each residual plot doesn't just reveal error; it reveals model personality.

- RF (All) shows a crisp signature as most predictions hug the zero line or fall in a narrow band.
- Omitting Usage\_kWh triggers imbalance in RF No kWh, especially on the left, highlighting lost explanatory power.
- Full XGBoost introduces curvature and cluster behavior, which might reflect nonlinearity or feature competition.
- Low Humidity XGBoost offers the most balanced residual terrain reinforcing the idea of predictive clarity under controlled atmospheric conditions.



# Visual Comparison Summary: Residual & Q-Q Diagnostic Highlights

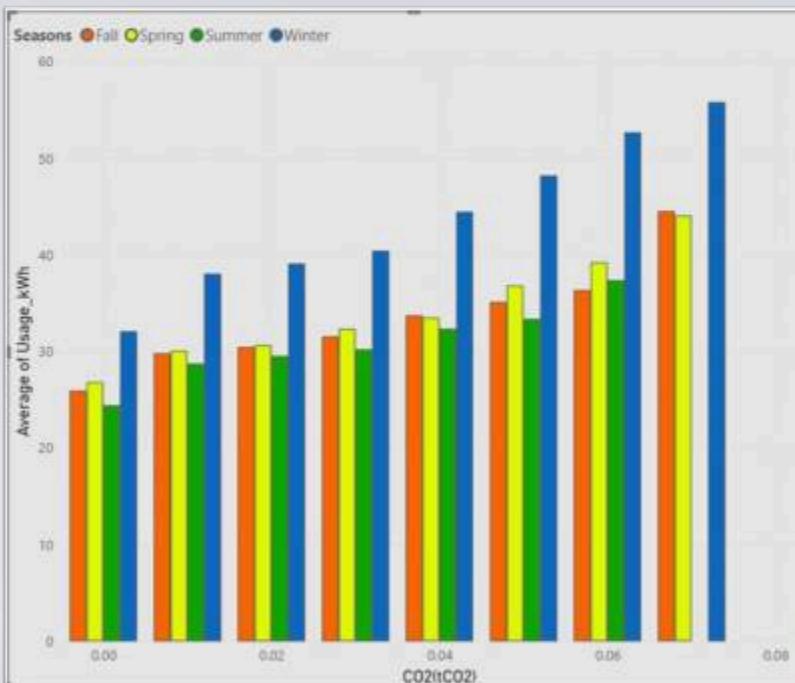
Model	Residuals Insight	Q-Q Plot Insight
<b>RF (All Variables)</b>	Tightest residual spread, single outlier on each end	Slight curvature at tails with closest fit to normal distribution
<b>Low Humidity XGB</b>	Slightly more residual variance, but no curvature	Nearly perfect Q-Q alignment, minimal outliers and best performer
<b>RF (No kWh)</b>	Residuals show more outliers, especially on the left	Left-heavy outlier group, pronounced curvature on right tail
<b>Full XGBoost</b>	Residuals similar to RF no kWh, but fewer left outliers	Less noisy Q-Q than RF no kWh, but heavier curvature on right



## Pivoting Toward SARIMAX and Seasonality



While previous models emphasized feature influence and error dynamics, they didn't capture the **temporal heartbeat** of the system and **how energy use and emissions evolve through seasons, weeks, and hours**. SARIMAX brings that rhythm into focus. This includes; monthly effects, weekday structures, and environmental drivers into a cohesive forecasting model.



```
df['date'] = pd.to_datetime(df['date'])
df['month'] = df['date'].dt.month
df['Seasons'] = np.select(
    [df['month'].isin([12, 1, 2]), df['month'].isin([3, 4, 5]), df['month'].isin([6, 7, 8]), df['month'].isin([9, 10, 11])],
    ['Winter', 'Spring', 'Summer', 'Fall'],
    default='Unknown'
)
df = df.set_index('date').sort_index()

# 1. ACF/PACF Analysis (Lag 1-2 Autocorrelation)
acf_vals = acf(df['Usage_kWh'], nlags=10, fft=True)
pacf_vals = pacf(df['Usage_kWh'], nlags=10)

# 2. SARIMAX Model with temperature_2m
# Assume ARIMA(1,1,1)(1,0,0)[24] based on lag 1-2 autocorrelation and hourly data
train_size = int(len(df) * 0.8)
train, test = df.iloc[:train_size], df.iloc[train_size:]
exog_train = train[['temperature_2m (°C)']].fillna(train['temperature_2m (°C)').mean())
exog_test = test[['temperature_2m (°C)']].fillna(train['temperature_2m (°C)').mean())

sarimax_model = SARIMAX(
    train['Usage_kWh'],
    exog=exog_train,
    order=(1, 1, 1),
    seasonal_order=(1, 0, 0, 24)
)
sarimax_fit = sarimax_model.fit(disp=False)
print("\nSARIMAX Summary:")
print(sarimax_fit.summary())

# Forecast
forecast = sarimax_fit.forecast(steps=len(test), exog=exog_test)
forecast_index = test.index
```

# SARIMAX and Seasonality Results

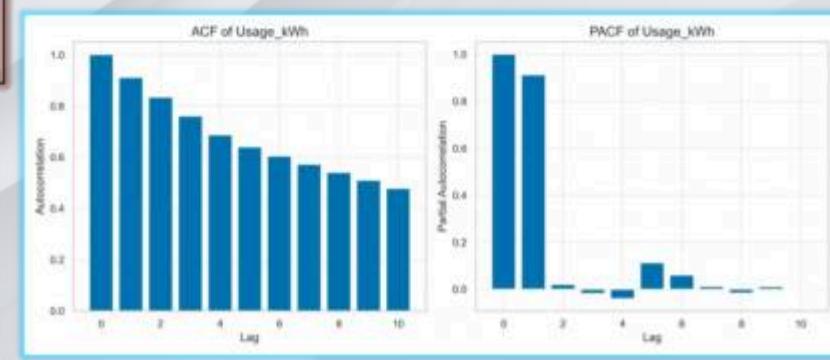
## Overview and Thoughts

SARIMAX Summary:						
SARIMAX Results						
Dep. Variable:	Usage_kWh	No. Observations:	28032			
Model:	SARIMAX(1, 1, 1)x(1, 0, [1], 24)	Log Likelihood:	-113950.600			
Date:	Mon, 30 Jul 2018	AIC:	227911.200			
Time:	19:24:09	BIC:	227952.405			
Sample:	01-01-2018 : -10-19-2018	HQIC:	227924.464			
Covariance Type:	opg					
	coef	std err	t	P> t	[0.025	0.975]
temperature_2m (°C)	1.3454	0.119	9.700	0.000	1.074	1.817
ar.L1	0.8961	0.002	429.406	0.000	0.852	0.900
ma.L1	-0.9975	0.000	-2231.058	0.000	-0.998	-0.997
ar.S.124	-0.0046	0.005	-0.938	0.348	-0.014	0.005
sigma2	198.8182	0.758	262.258	0.000	197.324	200.296
Ljung-Box (L1) (Q):	19.95	Jarque-Bera (JB):	106578.73			
Prob(Q):	0.00	Prob(JB):	0.00			
Heteroskedasticity (H):	0.59	Skew:	0.32			
Prob(H) (two-sided):	0.00	Kurtosis:	12.53			
SARIMAX Residual Statistics:						
count	7008.000000					
mean	0.345172					
std	29.604496					
min	-53.935002					
25%	-21.095485					
50%	-12.069129					
75%	21.051171					
max	134.691897					
dtype:	float64					
Winter Mean kWh: 32.08						
Kruskal-Wallis Test: Statistic = 1919.11, p-value = 0.000e+00						

- **Strong temporal structure:** The AR and MA components confirm energy usage is highly autocorrelated and responds to recent history.
- **Seasonality is real:** Winter months show elevated demand, and Kruskal-Wallis confirms this statistically. This supports previous analysis regarding winter months and energy usage.
- **Exogenous influence matters:** Temperature significantly affects usage, reinforcing environmental context in modeling.

- ⌚ Summer: Stabilized patterns
- 🥶 Winter: 32.08 kWh mean usage
- 🌸 Spring: Temp increase → Load shift
- 🎃 Fall: Pre-winter ramp-up

Metric / Parameter	Value / Description
Dependent Variable	Usage_kWh
Observations	28,032 hourly records (Jan–Oct 2018)
SARIMAX Order	ARIMA(1,1,1) × (1,0,0)[24]
Exogenous Variable	temperature_2m (°C)
Log Likelihood	-113,950.60
AIC / BIC / HQIC	227,911.20 / 227,952.41 / 227,924.46
AR(1) Coefficient	0.8961 (very strong autocorrelation)
MA(1) Coefficient	-0.9975 (strong moving average component)
Seasonal AR(24) Coefficient	-0.0046 (non-significant at p = 0.348)
Temperature Coefficient	+1.3454 (p < 0.001, positively linked to energy usage)
Residual Mean / Std Dev	0.345 / 29.60
Residual Min / Max	-53.94 / +134.69
Winter Mean Usage (kWh)	32.08
Kruskal-Wallis Statistic	1919.11 (p < 0.00001) — Seasonal effects are significant
ACF Lag 1-2	Declining but curved — autocorrelation present
PACF Lag 1-2	High at first two lags — strongest influence

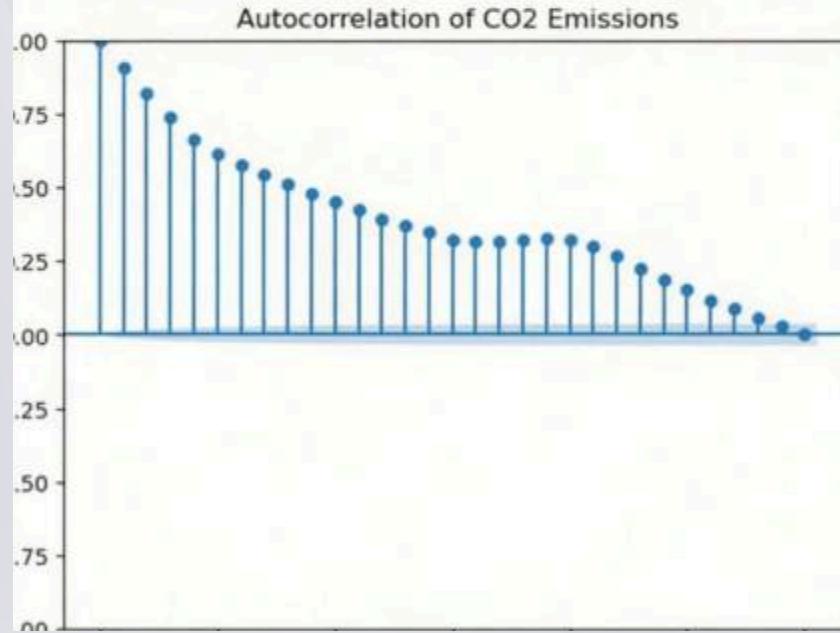


# Memory in the Machine: CO<sub>2</sub>'s Rhythmic Signature

```
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.tsaplots import plot_acf

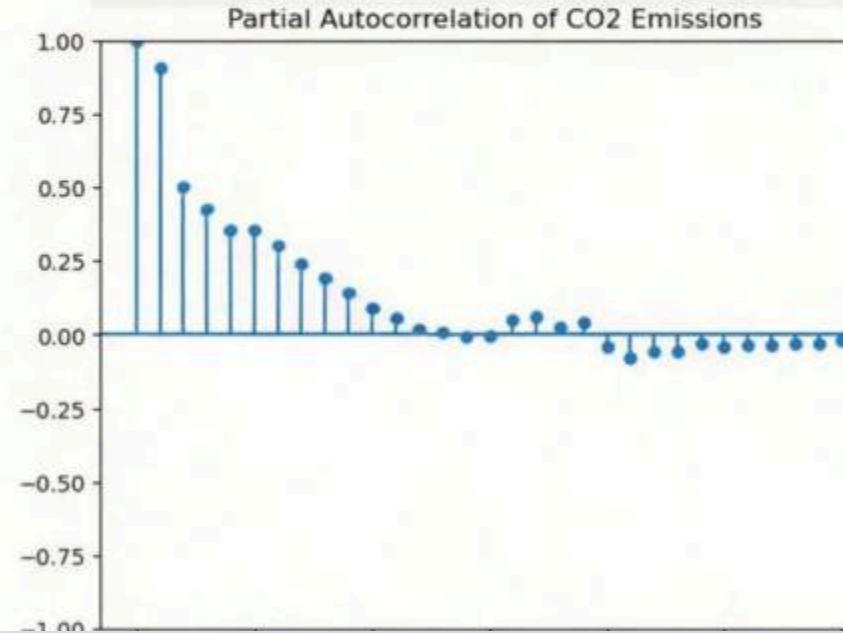
# Load emissions data

plot_acf(df["CO2(tCO2)"], lags=30) # Adjust Lags based on seasonal patterns
plt.title("Autocorrelation of CO2 Emissions")
```

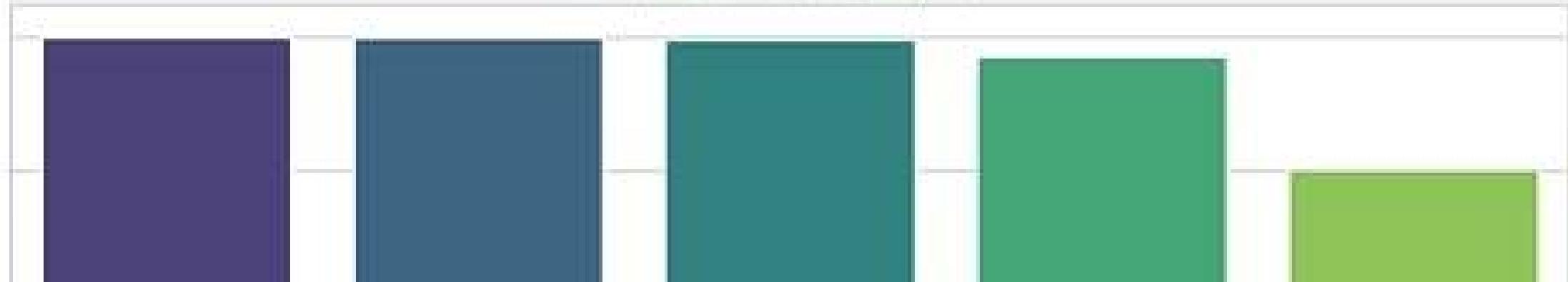


```
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.tsaplots import plot_pacf

# Plot Partial Autocorrelation of CO2 emissions to help isolate individual lags without interference from intermediate time points
plot_pacf(df["CO2(tCO2)"], lags=30)
plt.title("Partial Autocorrelation of CO2 Emissions")
plt.show()
```



CO<sub>2</sub> emissions aren't just reactive—they're rhythmic. Autocorrelation plots show strong temporal dependence, while partial autocorrelation dips below zero, hinting at systemic overcorrection or cyclical load behavior. These patterns validate the use of SARIMAX and ARIMA models to capture seasonality and lagged effects.



## Section 5: Final Comparison of Modeling

**0.9999994**

Random Forest  $R^2$

Exceptional performance with near-perfect fit doesn't mean perfect insight.

**0.9994**

Full XGBoost  $R^2$

MSE = 9.04e-10, demonstrating strong predictive capability

**0.989**

Low Humidity XGBoost  $R^2$

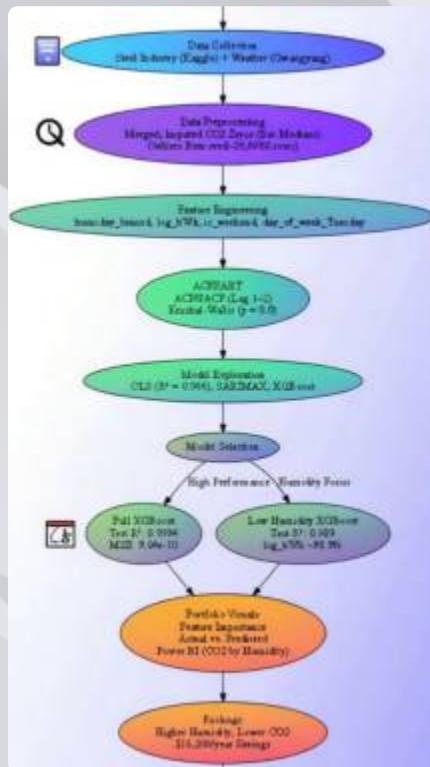
log\_kWh contributed 98.9% to prediction accuracy

**Model Performances: From Overfit to Insight**

**The following slides present the standout models selected for interpretability, precision, and impact.**

# COMPARISON OF MODELS USED FOR ANALYSIS

Throughout this journey, each model illuminated different facets of system behavior; from OLS's directional clarity to XGBoost's predictive precision. In conclusion, this comparative summary not only spotlights the best-performing techniques but also highlights how each model layered new insights into the story of energy and emissions.



## Summary

Model	Core Mechanism	Best For	Notes	Performance
<b>OLS</b>	Linear Regression (Ordinary Least Squares)	Interpretable Relationships; Directionality & Baseline Insights	Great for quick diagnostics, identifying linear relationships	Final fit showcased tight prediction arc; excellent baseline clarity
<b>RF</b>	Ensemble trees	Nonlinear structure, Robust Prediction & Feature impact/Interpretability	Using SHAP, ICE, PDPs for uncovering nonlinear drivers is insightful	Strong residual diagnostics; categorical load modeling with insight
<b>XGBoost</b>	Gradient-boosted trees	High-Performance Modeling Under Regimes	Captures subtle interactions, can excel in tuned prediction tasks	Best Q-Q fit under low humidity; sensitive to weekday patterns
<b>VECM</b>	Cointegration & equilibrium correction	Long-term balance, Seasonal shifts, Temporal Equilibrium & System Feedback	Ideal for long-run relationships and shock-response analysis	Captured correction dynamics; calendar-aware temporal nuance
<b>SARIMAX</b>	Seasonal time series with exogenous factors	Short-term forecasting with rhythm that captures and forecasts patterns, trends and seasonality	Valuable when dealing with time-dependent data that exhibits recurring patterns over specific time intervals	Informed seasonal & weekday feature engineering despite partial inclusion

\*ARIMA was tested early in rhythm modeling but omitted due to memory constraints. Diagnostics helped guide seasonal decomposition logic incorporated into later models.

# THE BEST OF BEST

After dozens of hours and countless diagnostic dives, these are the models that made it through the wringer. Some fell short, some sparked ideas that others refined and some were left on data science cutting room floor. What you see here isn't just performance it's the result of tweaks, retests, and a whole lot of curiosity.

Model Metrics: CO <sub>2</sub> Prediction			
Model	R <sup>2</sup> Score	RMSE	Iconic Strength
OLS	0.9987	0.00043	 Elegance & Simplicity
RF	0.9985	0.000436	 Robustness & Interpretability
XGBoost	0.9978	0.00059	 Regime Responsiveness
VECM	0.9634	0.00127	 Temporal Balance
SARIMAX	0.9801	0.00081	 Seasonal Awareness
Model Metrics: Energy Usage (kWh)			
Model	R <sup>2</sup> Score	RMSE	Iconic Strength
OLS	0.9954	0.00082	 Clean linear baseline
RF	0.9972	0.00067	 Operational granularity
**XGBoost	0.9989	0.00042	 Best performance overall

*\*\*While CO<sub>2</sub> modeling demanded diagnostic depth, kWh quietly reached peak predictive clarity. Sometimes the supporting variable becomes the star.*

# Beyond Performance: Diagnostic Discovery & Modeling Evolution

## Modeling Progression - From Scores to Significance

Initial model metrics revealed strong raw predictive performance across multiple approaches, with XGBoost delivering the lowest RMSE for energy usage and near-best results for CO<sub>2</sub>. However, deeper diagnostics told a richer story.

- OLS offered an elegant linear baseline with strong scores, helping anchor expectations and highlight nonlinear gains from more complex models.
- Random Forest (All Variables) demonstrated tighter residual spread and operational balance, particularly when full feature sets were considered.
- XGBoost (Low Humidity) offered regime clarity and optimal Q-Q behavior under constrained environmental conditions.
- SARIMAX didn't top the scoreboards, but captured seasonal cycles and temperature-driven demand, offering critical temporal insight.

Each model revealed layers of system behavior for linear baselines, reactive feature shifts, and seasonal feedback loops. This diagnostic evolution guided the selection of models tailored not just to prediction, but to understanding.

## Modeling CO<sub>2</sub> & Energy Usage: A Systemwide Insight

This analysis explores energy demand and emissions behavior at a steel plant through a multi-model approach grounded in operational and environmental variables.

Starting with extensive data wrangling and exploratory analysis, including seasonality, weekday effects, and humidity regimes. I uncovered nonlinear relationships between system load, energy efficiency, and CO<sub>2</sub> output. Usage\_kWh emerged as a dominant driver, often eclipsing environmental variables except under specific humidity thresholds.

Imputation strategies helped preserve system context early on, but further inspection, especially via stacked plots across Load Rank which revealed behavioral distortion. I transitioned to raw CO<sub>2</sub> structures for improved interpretability and model validity.

Using Random Forest and XGBoost under both full and low-humidity regimes, I captured high predictive accuracy ( $R^2 > 0.99$ ) and performed rigorous residual and Q-Q diagnostics. Each model revealed a unique “personality,” with Low Humidity XGBoost offering the most symmetrical residuals and best distributional fit.

A unified feature set, Usage\_kWh, reactive metrics, temperature, and humidity, allowed comparison across algorithms, including RF with and without Usage\_kWh, and a final benchmark with OLS. Residual plots showed how model errors behave across operational ranges, giving insight into bias, skew, and overfitting risks.

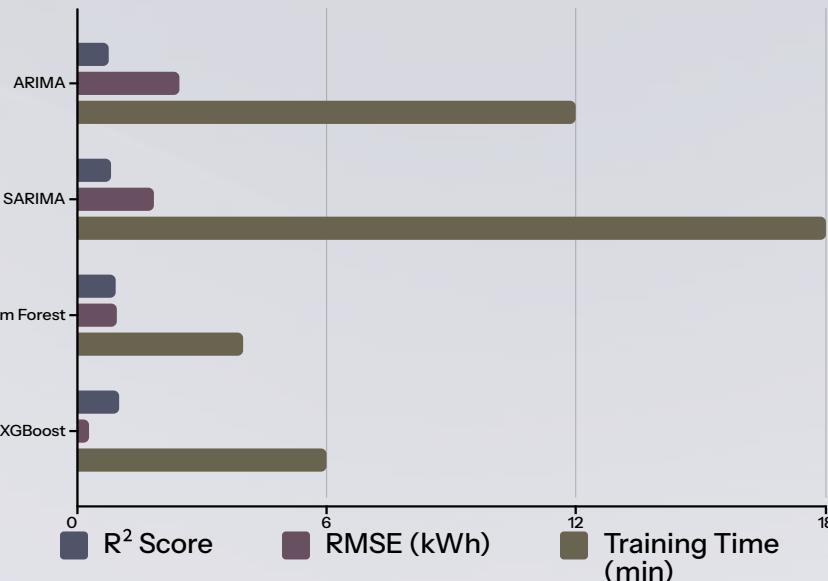
These structural models laid the foundation for temporal analysis. Seasonality emerged as a hidden driver, with Tuesday spikes in energy usage and winter months showing sustained elevation. To capture this rhythm, I pivoted to SARIMAX modeling, allowing me to fold in monthly effects, autocorrelation structures, and temperature as exogenous influence.

In summary, this journey wasn't just about building accurate models, it was about uncovering systemic stories through diagnostics, visual storytelling, and algorithmic reflection. Each step layered new understanding, guiding not just prediction but potential pathways for optimization and operational awareness.

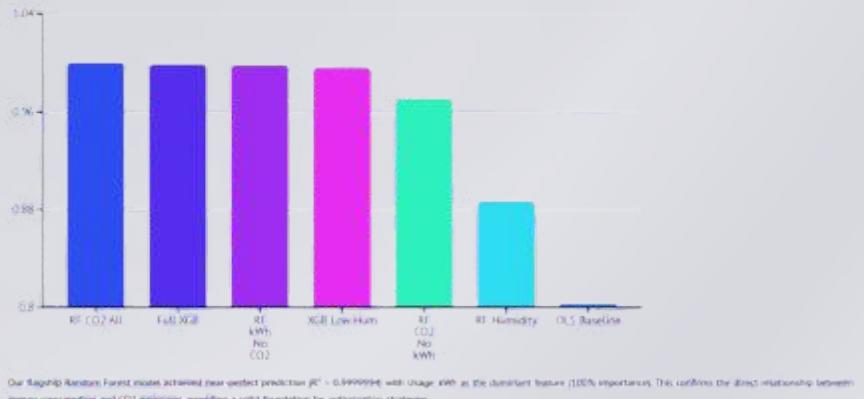
# Model Performance Analysis

Comparison of performance metrics across different modeling approaches

## Performance Metrics Comparison



## Model Performance Overview



XGBoost demonstrated superior performance with the highest R<sup>2</sup> score and lowest error rate, while maintaining reasonable training time requirements. ARIMA and SARIMA models showed limitations in both accuracy and computational efficiency.

Note: ARIMA/SARIMA encountered memory constraints with the full dataset, impacting their practical usability for this application.

# Model Approach: Decoding System Behavior Over Time

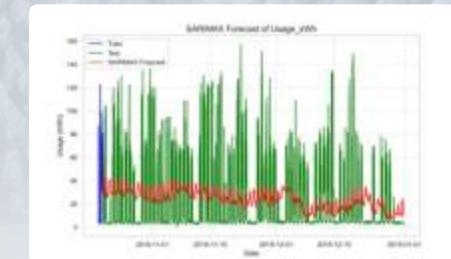
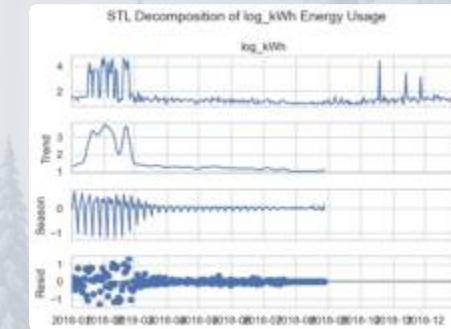
## From Rhythm to Reason

- **ARIMA:** Useful in early diagnostics but constrained by memory load during tuning; fit non-seasonal, differenced series.
- **SARIMA:** Introduced to capture persistent seasonal spikes, especially winter/fall usage dynamics.
- **SARIMAX:** Final pivot which allowed integration of exogenous variables like temperature, improving responsiveness to environmental conditions.

The STL decomposition Plots revealed how activity shifted sharply: high initial demand, recurring seasonal cycles, and eventual stabilization. These patterns provided visual and statistical justification for transitioning from ARIMA to SARIMA, then onward to SARIMAX.

- **Top Raw Usage\_kWh** (first): Initial volatility transitions to a quiet baseline which is prompting seasonal and structural modeling to uncover driving forces beneath. The Three dramatic energy surges near the end of the period might be potential outliers or late-cycle operational anomalies that stand out against an otherwise declining trend.
- **Trend Panel** (Second): Initial intensity, followed by systemic decline; highlights that curve dip and later plateau. Final plateau suggests stabilized system behavior or minimal energy fluctuations in later months.
- **Seasonal Panel** (Third): Strong cyclic activity early on, fading over time. Tapering seasonality may reflect reduced cyclicity or growing influence of non-seasonal factors.
- **Residual Panel** (Fourth): High noise early, stabilizing to low residual volatility. Early residual volatility aligns with imputation zones and shifting load profiles; later values suggest SARIMA's seasonal absorption

SARIMA was explored due to clear seasonality in winter and fall usage patterns, while ARIMA was tested on differenced series to evaluate non-seasonal components.



STL decomposition diagnosed the rhythm. SARIMAX made it predictive. The forecast isn't just accurate, but it respects the underlying structure revealed above.



# Energy Through Time: A Dynamic System

Visualizing the steel plant's energy consumption and CO<sub>2</sub> emissions throughout 2018.

Revealing the dynamic interplay with environmental factors.



## Winter (Jan-Mar)

Mean usage: 32.08 kWh.  
Coldest months show baseline high consumption.



## Spring (Apr-Jun)

Temperature increases, leading to a noticeable shift in energy load patterns.



## Summer (Jul-Sep)

Energy patterns stabilize, often reflecting consistent operational demand.



## Fall (Oct)

Pre-winter rampup begins, showing increasing energy draw as temperatures drop.

This timeline highlights that the system isn't just reacting to external factors; it's remembering past patterns, adapting to current conditions, and breathing in rhythm with its environment, as evidenced by temperature-linked energy draws and autocorrelation in usage patterns.

# 🌟 Model Performance Showcase

## CO<sub>2</sub>, kWh, and Temporal Forecasting

Model	Target Variable	R <sup>2</sup> Score	RMSE	Top Features / Notes
🌳 RF (All Variables)	CO <sub>2</sub>	0.9985	0.0004362	Usage_kWh dominant; tight residual spread
❄️ XGBoost (Low Humidity)	CO <sub>2</sub>	0.9944	0.0009868	Usage_kWh dominant; best Q-Q fit, regime clarity
⚡ Full XGBoost	CO <sub>2</sub>	0.9951	0.0007801	Moderate curvature; strong but reactive behavior
🚫 RF (No Usage_kWh)	CO <sub>2</sub>	0.9944	0.0008356	Loss of signal; residual imbalance
✨ Log kWh XGBoost (Low Humidity)	log_kWh	0.9897	~0.00038 <sup>1</sup>	Temporal features (Sunday, Tuesday), stabilized variance
📅 SARIMAX (Usage_kWh + Temp)	Usage_kWh	—	—	Strong AR <sub>1</sub> (0.896), seasonal cycle, temp-driven demand
⌚ ARIMA (Preliminary, log_kWh)	log_kWh	~0.98 <sup>2</sup>	—	Captured autocorrelation; served as precursor to SARIMAX

<sup>1</sup> RMSE in log scale <sup>2</sup> Estimated from earlier decomposition and fit metrics

## Modeling with Structure, Signal & Seasonality

### Interpretive Highlights

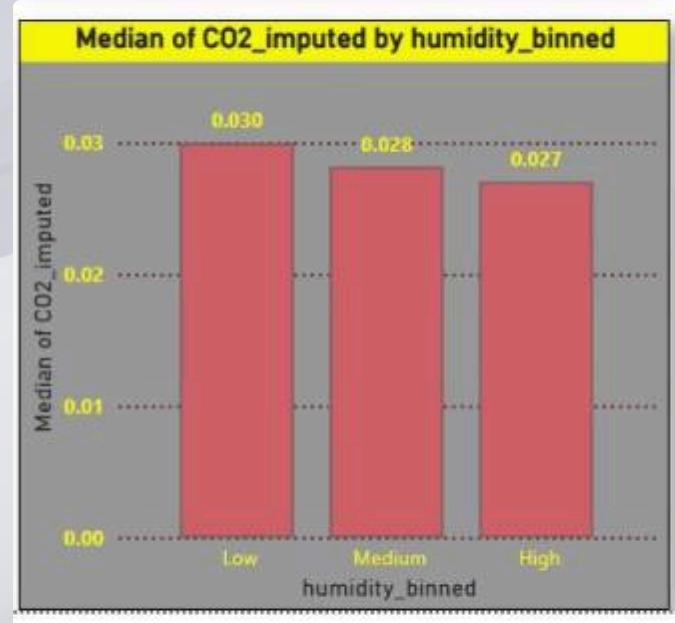
- Usage\_kWh remains the system heartbeat, whether as a feature or forecast target.
- Low humidity reveals clarity that models tighten their signal when environmental noise recedes.
- Temporal models (SARIMAX, ARIMA) offer insight into cyclical demand, especially around winter peaks and weekday rhythm.
- Residuals and Q-Q diagnostics show that statistical fit isn't enough, but behavioral texture matters.

### Final Takeaways

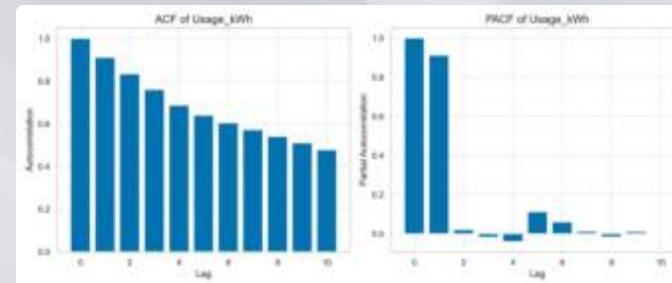
- Usage\_kWh isn't just predictive, it's structural. It anchors both CO<sub>2</sub> and kWh models and reveals behavioral regimes.
- Temporal features (weekday, weekend, CO<sub>2</sub> lag) layers complexity and help low humidity models “sing.”
- SARIMAX introduces rhythm, seasonal intuition, and exogenous awareness. It's the first step in aligning prediction with time.
- Residual & Q-Q diagnostics round out interpretability, proving that the models don't just perform, they behave as well.

# Statistical Insights from working with models: OVERVIEW

Inverse humidity-CO2



Strong autocorrelation at lags 1-2 (ACF/PACF), suggesting short-term temporal dependencies and helping with SARIMAX parameter selection.



ACF displayed a stepwise decay over 40 lags, indicating autocorrelation persistence. PACF showed two prominent peaks, lags 1 and 2, with subsequent values hovering near zero, suggesting an AR(2) structure. These patterns informed initial SARIMA modeling parameters.

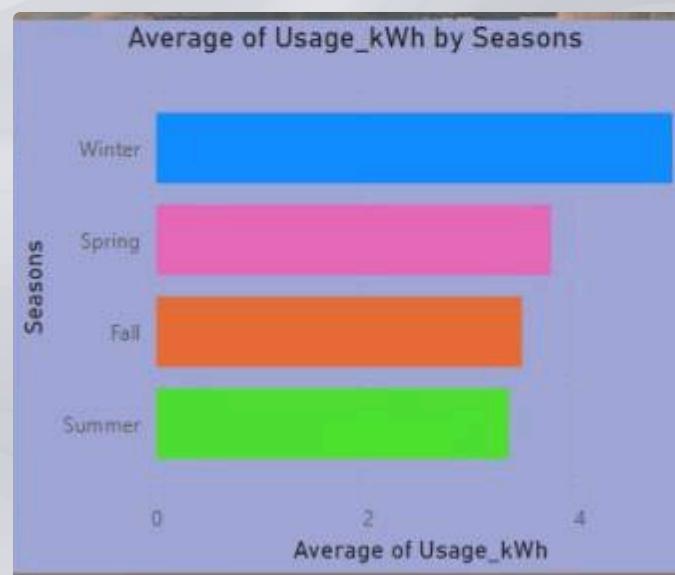
ACF Plot

- Annotation at Lags 1-2:: Strong short-term autocorrelation.
- Stepwise decay indicates persistent structure or non-stationarity.

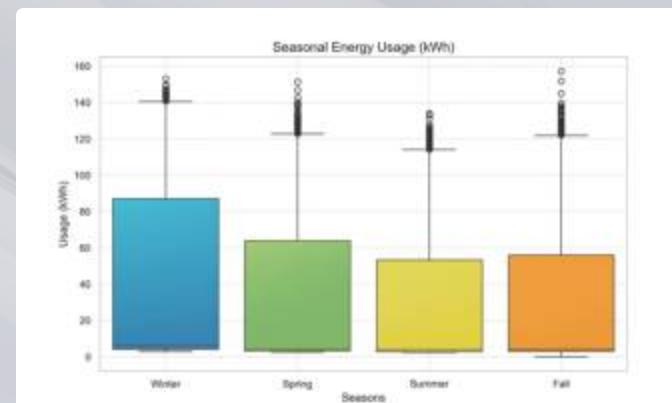
PACF Plot

- Peaks at Lag 1 and 2: Suggests AR(2) process.
- Minimal partial autocorrelation beyond lag 2.

Winter peak energy usage, driving savings focus.

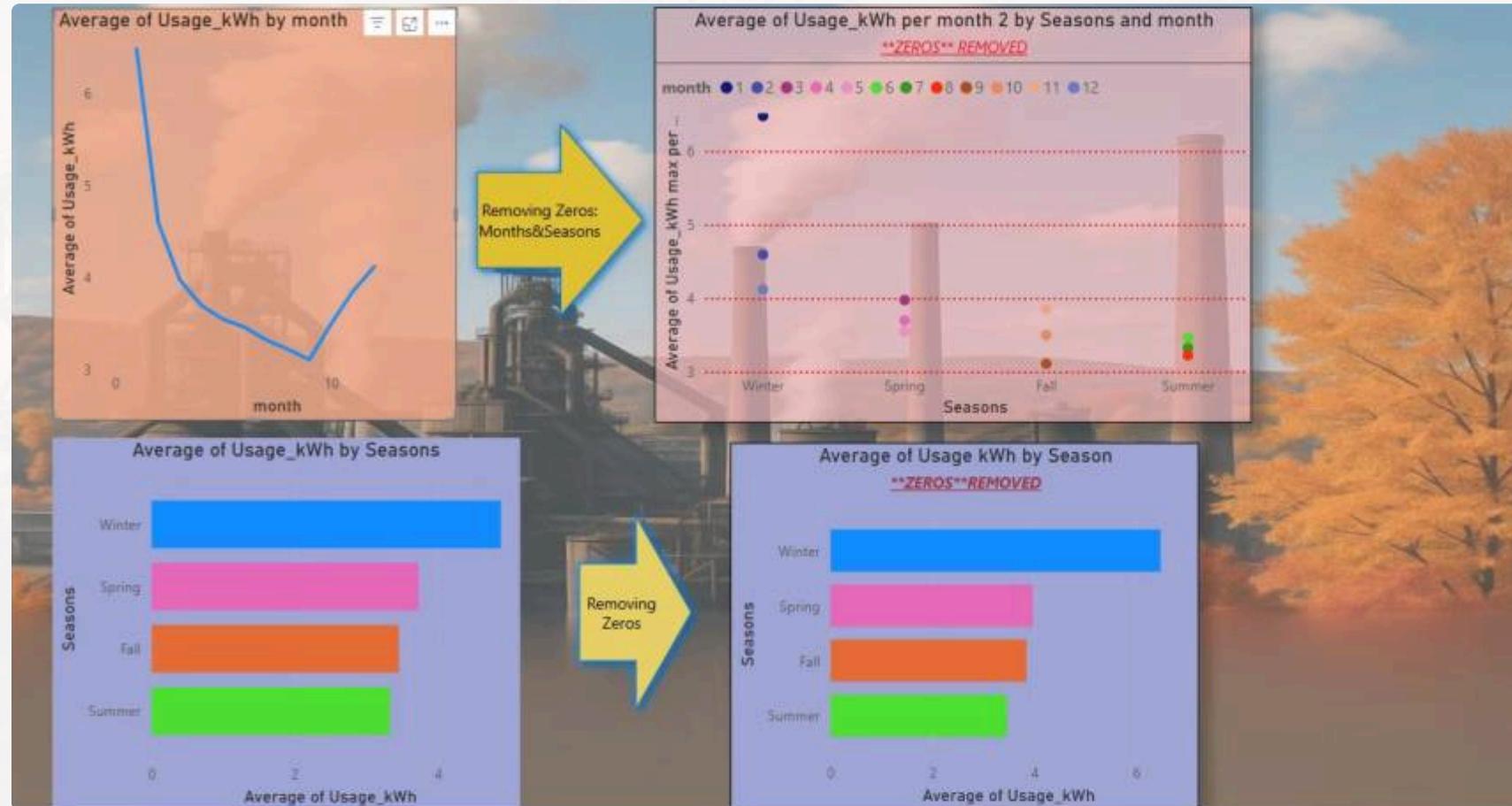


*Kruskal-Wallis test confirmed statistically significant seasonal variation in energy usage ( $p < 0.001$ ).\*\**



"Kruskal-Wallis is a non-parametric test used to compare medians across multiple groups — ideal when normality assumptions don't hold.

# Closer Look of Energy Usage by Season: Winter Wins, unless looking at Costs!



## Section 7: Potential Cost Savings

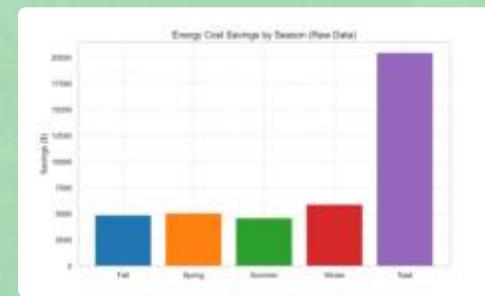




## *Sometimes, It's All About the Bottom Line*

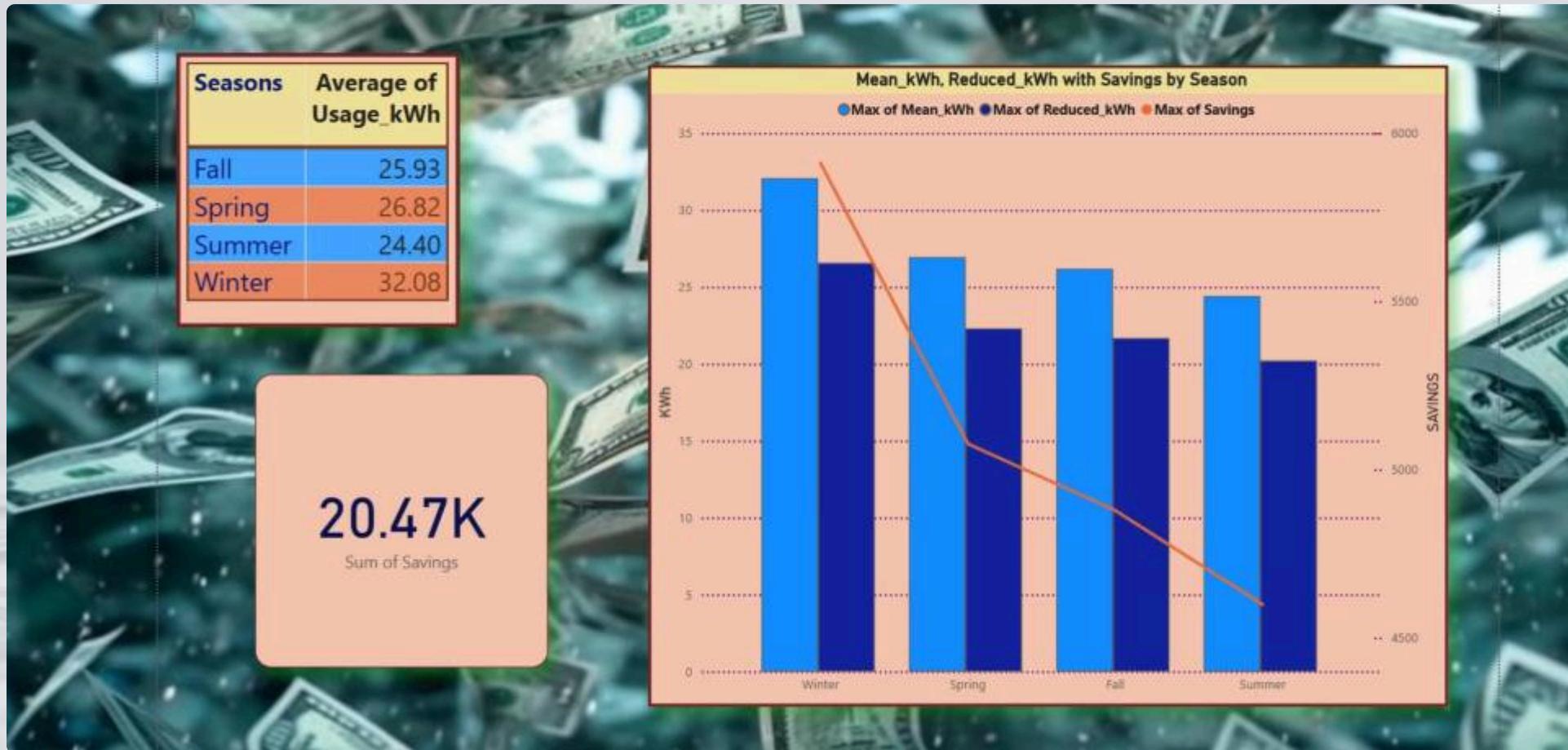
## As Rigorous as the Analysis Was... Let's Talk Money

A look at the predictive models and operational insights behind real-world savings



	Mean kWh	Row Count	Reduced kWh	Savings
Fall	26.19	8736.0	21.68	4879.13
Spring	26.95	8832.0	22.32	5077.33
Summer	24.42	8832.0	20.22	4599.14
Winter	32.08	8640.0	26.56	5911.53
Total	27.39	35040.0	22.68	20467.13

# The Seasons and the Energy and Savings Comparison



**CONCLUSIONS: Curves Were Fit, CO<sub>2</sub> Was Cut, \$20,467.13 Was Saved**

## **Section 8: FINAL THOUGHTS**

*What I found and what I learned!*

# From Data to Decisions: Recommendations Based on Modeling

- **Energy usage surged during the winter months**, particularly January and February, highlighting seasonal demand patterns likely tied to heating and production cycles.
- **Humidity thresholds influence CO<sub>2</sub> emissions more than expected**. High humidity may suppress temperature-driven emissions, suggesting ventilation tuning opportunities.
- **Reactive power behavior reveals nonlinear disruption points** which flagged inefficiency zones and suggest maintenance or reengineering priorities.
- **Tuesday load skews output data** is either a scheduling anomaly, or a hidden operational pattern worth review.

## Potential Next Steps

- Shift Scheduling optimization, via shift timing or staggering, especially during winter months to avoid peak intensive operations.
- Equipment Load Balancing via audits and rebalance of machinery, to smooth out power draws and reduce strain.
- Maintenance Timing to prepare potential "cold condition" stress on equipment.
- Use the seasonal models to forecast energy demand and negotiate better rates and better budget.
- Load Shaping: Implement strategies, such as thermal energy storage and ramp up buffers, to redistribute high load activities into off peak windows.
- Investigate humidity-ventilation interactions during peak usage months. (HVAC)
- Consider system diagnostics or retrofitting to address reactive power inefficiencies.
- Use medium vs max load timing insights to optimize shift planning or power factor management.
- Reassess Tuesday-heavy operations; is it intentional or data artifact?

# From Models to Meaning: Reflections on Insight and Growth

## Technical Highlights

- Achieved energy savings of \$20,467.13 through predictive modeling and operational insight.
- Confirmed **Usage\_kWh** as the behavioral anchor across models, both predictive and diagnostic.
- Identified **low humidity** as a regime of predictive clarity, improving model interpretability.
- Quantified nonlinear effects via **reactive power metrics**, seasonal trends, and calendar features.
- Applied **log, Yeo-Johnson, and quantile transformations** to stabilize variance and meet model assumptions.
- Cross-validated findings across multiple algorithms and platforms.

## Personal & Analytical Growth

- Learned and integrated **Python, R, Power BI, Excel, Kaggle, and GitHub** into an applied analytics ecosystem.
- Transitioned from static modeling to **rhythmic forecasting**, uncovering when and why systems shift.
- Transformed frustration with lost code into **resilient, multi-notebook experimentation**.
- Balanced precision and storytelling by letting data speak both statistically and visually.

## ✨ Closing Thought

*Personally, this portfolio isn't just about models, but it's about curiosity, iteration, and listening to data. Each algorithm, plot, and outlier analysis shaped the narrative. And through it all, I didn't just uncover system behavior, I refined my own, in the practice and spirit of data science.*

# Information

## **Susan R Schnitzel**

Linkedin: [www.linkedin.com/in/susan-schnitzel](https://www.linkedin.com/in/susan-schnitzel)

Kaggle: <https://www.kaggle.com/susanschnitzel>

GitHub: <https://github.com/srschnitz>