

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 CO₂ 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



Data Collection & Transformation

Kaggle dataset + weather merge, datetime alignment, CO₂ 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

Sections at a Glance



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

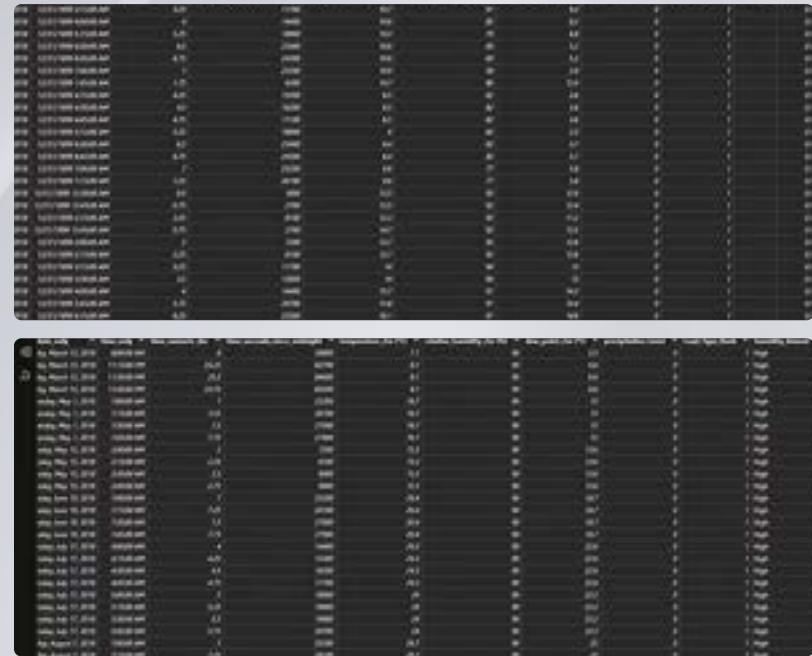
DATA WRANGLING

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.



The image shows two side-by-side screenshots of a data processing or visualization tool. Both screens display a grid of data rows. The top screen has columns labeled with various parameters such as date, day, month, and year. The bottom screen shows a similar grid but with more columns, likely representing additional variables or transformed data. The data appears to be time-series information, possibly related to energy consumption and weather.

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

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

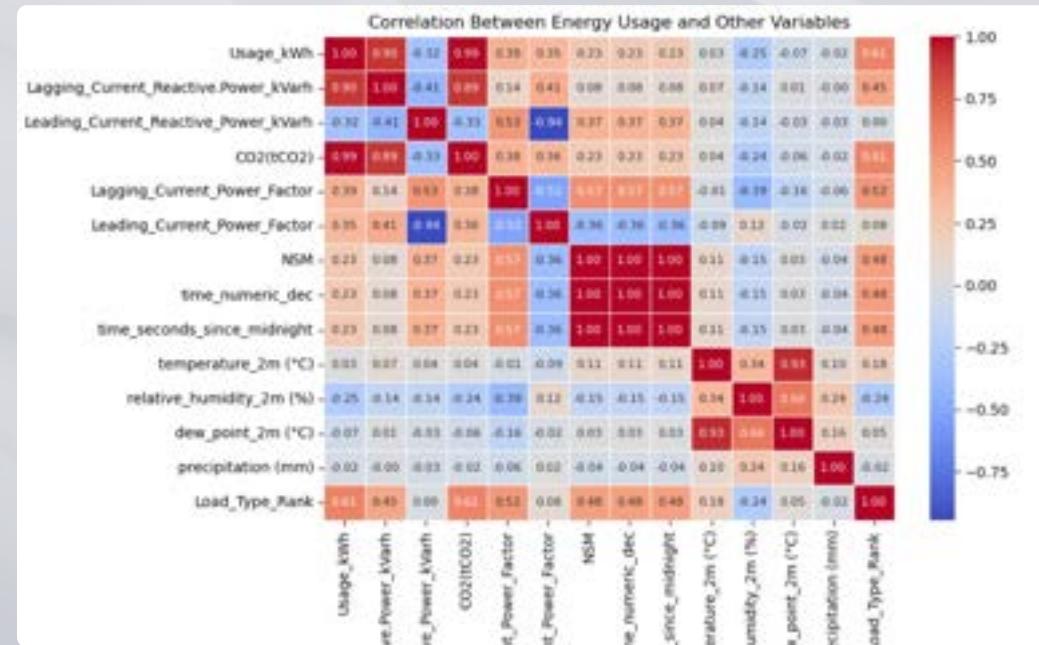
This slide may be brief, but the work behind it wasn't. Data merging, scaling, and transformation took plenty of sweat, tweaks, and caffeinated perseverance; none of which fit neatly in a bullet point.

Section 2: Exploratory Analysis and Insights

Exploratory and Data Insights

Getting to know the Data

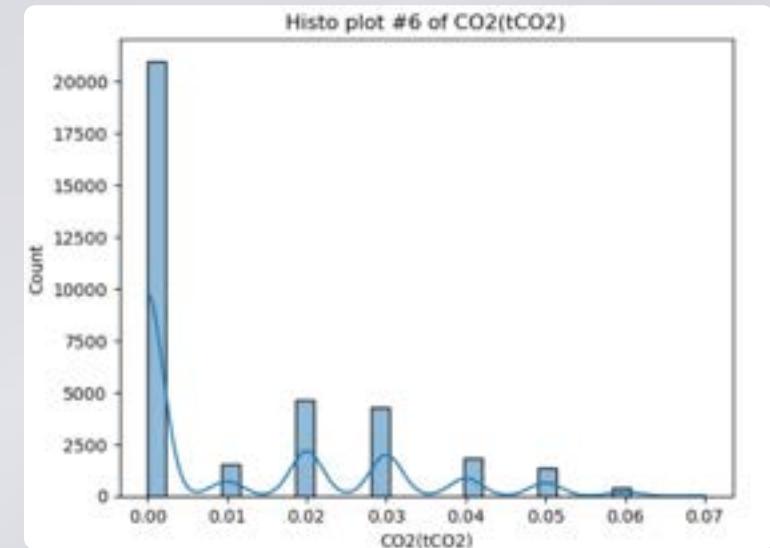
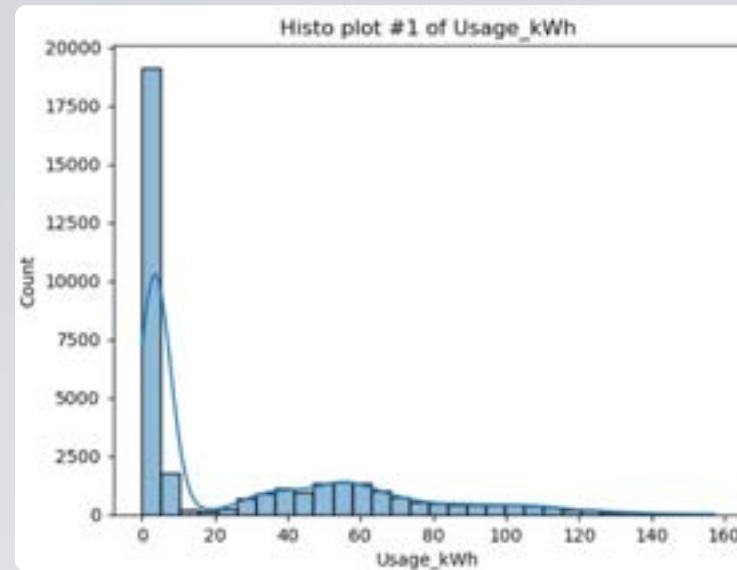
Before diving into modeling, I always take time to explore and interact with the data, not just to clean it, but to understand it. This includes visualizing distributions, examining trends over time, and checking for correlations or unexpected patterns. These early insights often spark better feature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



- Initial heatmap analysis revealed a strong positive correlation between energy usage (kWh) and CO₂ emissions; a relationship consistent with expected physical output, where increased energy production corresponds with higher emissions. This was key to guiding modeling decisions.

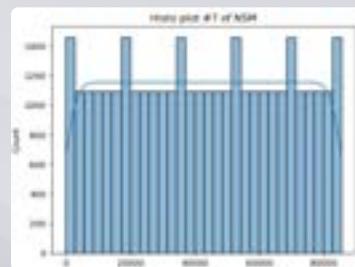
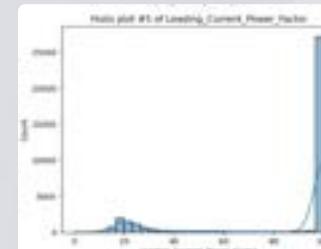
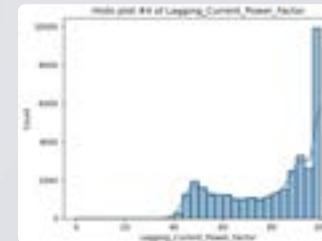
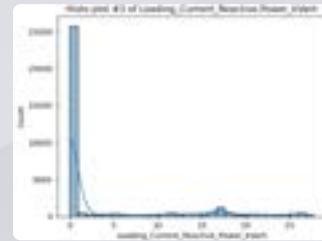
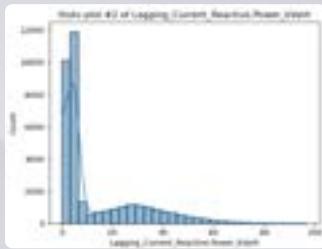
Exploratory Data Insights - Initial Distributions

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.

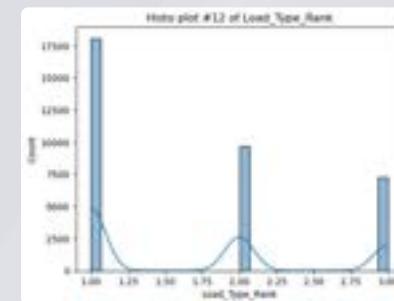


In this case, histograms of CO₂ and kWh revealed skewed distributions, leading to the application of transformation (log) on energy usage and an imputation strategy for zero-inflated CO₂ values.

Exploratory Data Insights: Energy Distributions

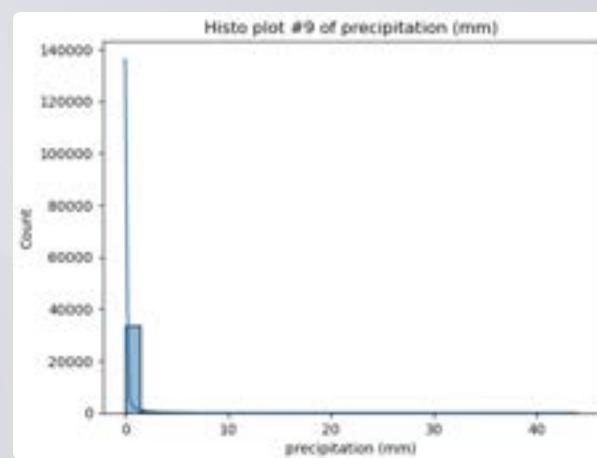
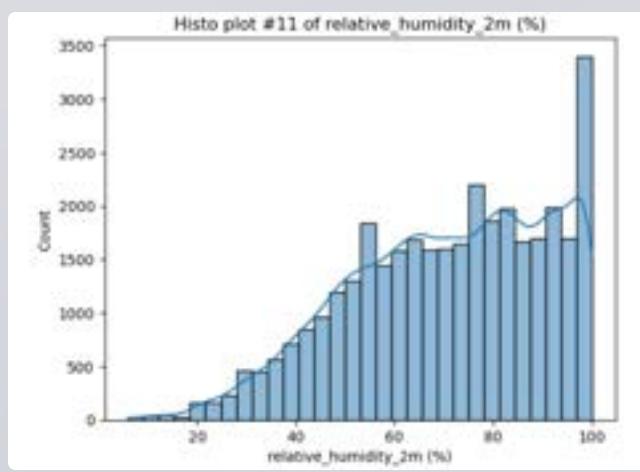
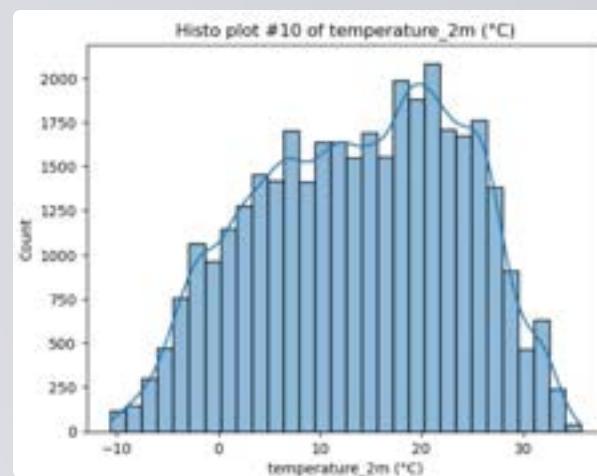
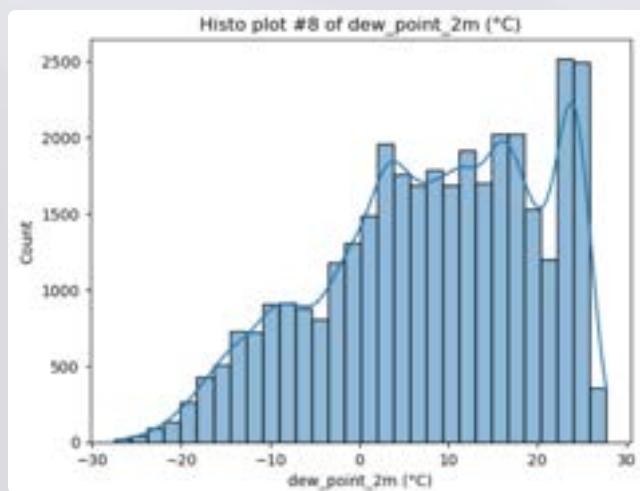


While several energy-related variables exhibited skewed distributions with some leaning heavily toward higher values. I also explored features like NSM and Rank, which are inherently nonparametric in nature. > > As someone with a background in statistics, I find it important not only to identify when 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

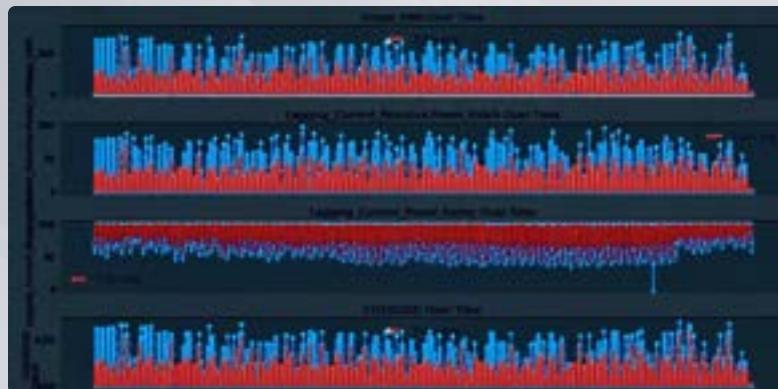


Histogram analysis showed that temperature and dew point followed approximately normal distributions, supporting their direct use in modeling without transformation. In contrast, humidity exhibited a strong right skew, suggesting the site regularly experiences high-humidity conditions. This insight guided 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.

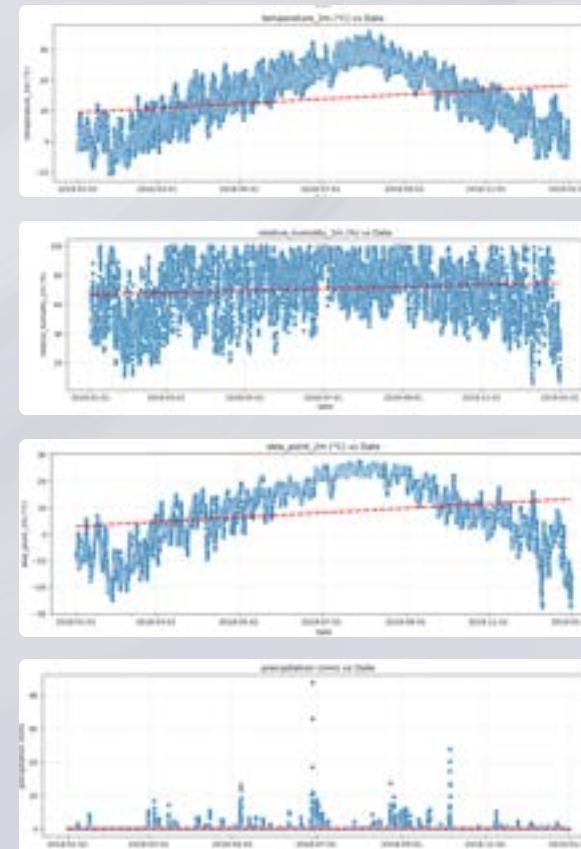
TIME SERIES OF VARIOUS VARIABLE

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

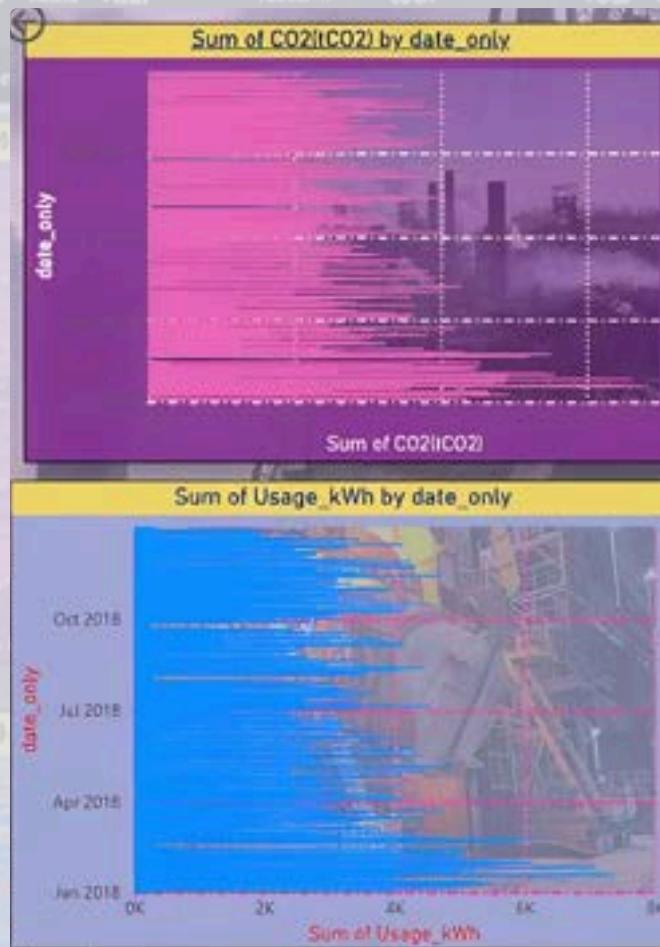
Energy and CO2 variables



Environmental /Weather variables

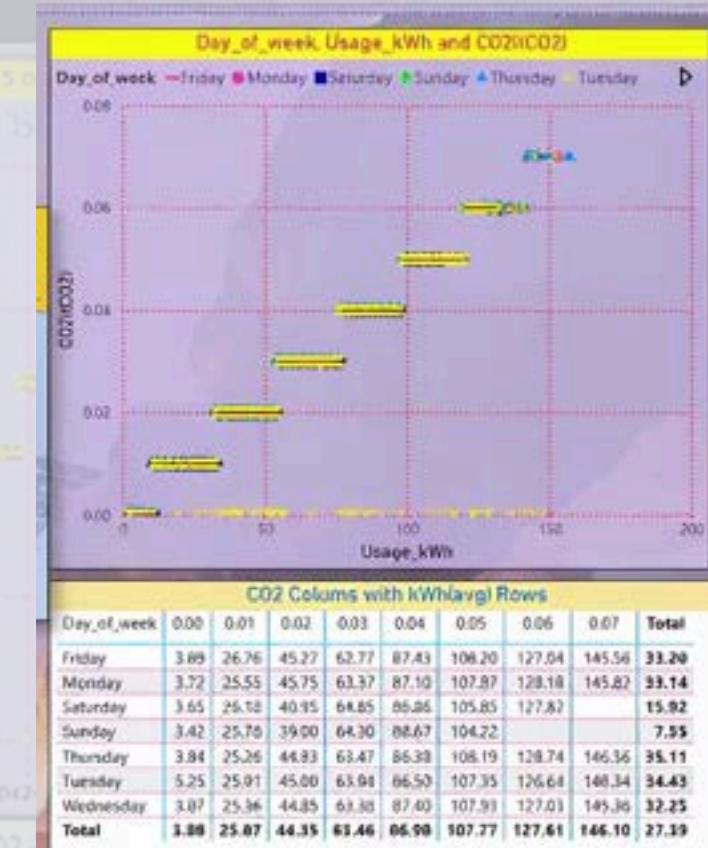


Time Series Comparison of CO2 and Energy Usage kWh



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..

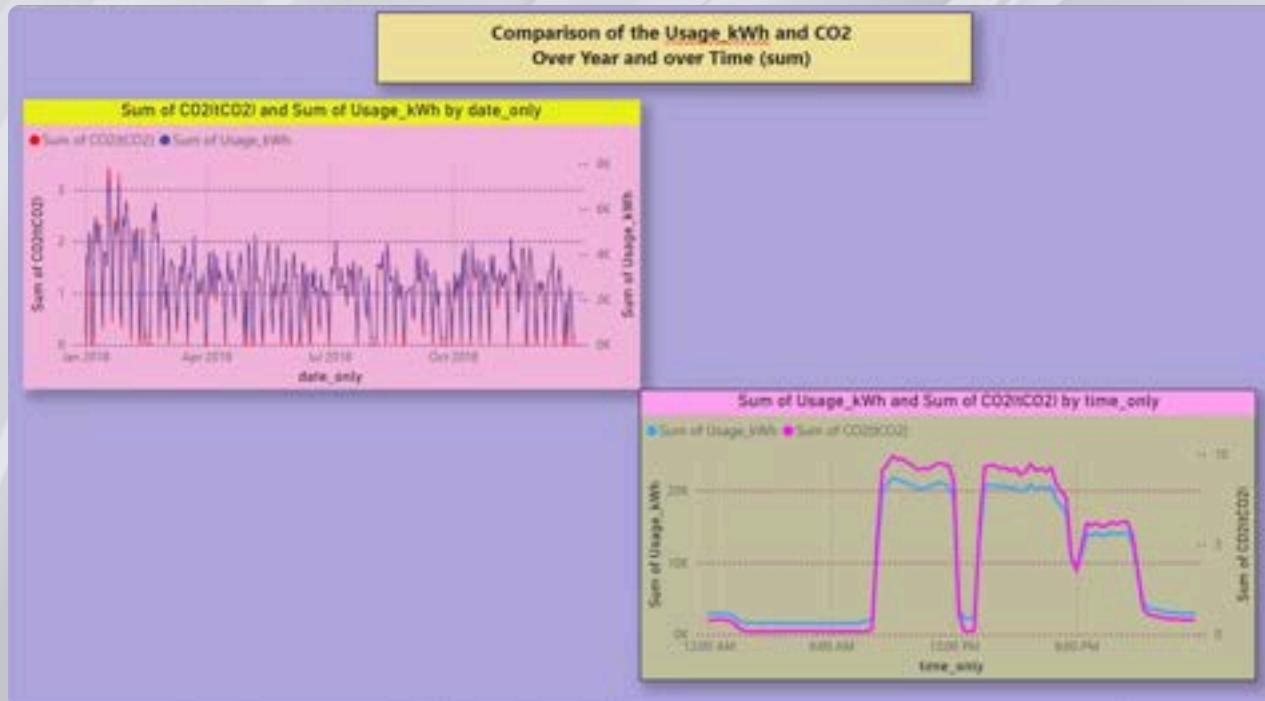


Days of The Week

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.”

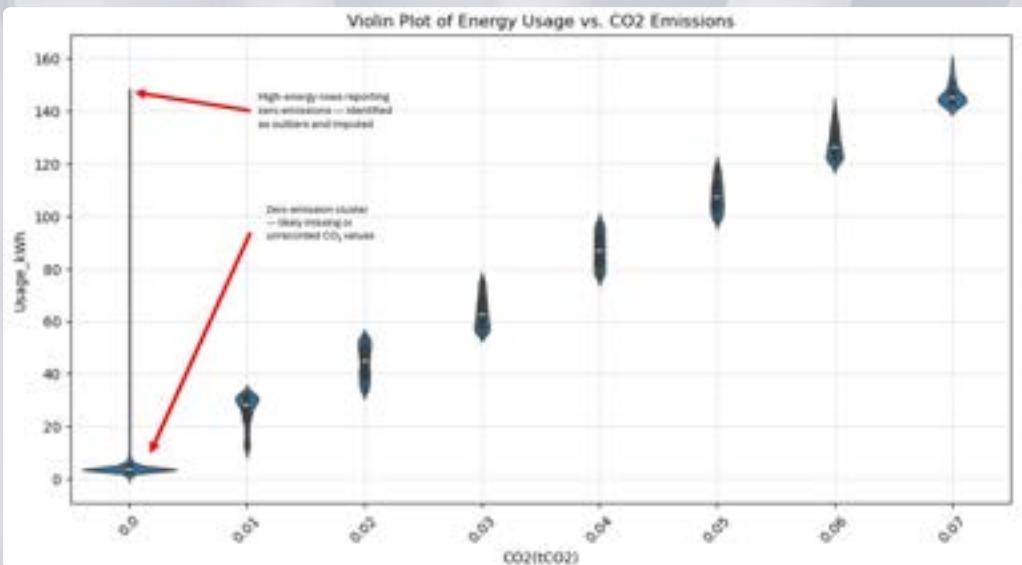
This correlation aligns with thermodynamic expectations – greater energy conversion yields higher CO₂ output. Also, it shows a large amount of zeros for the CO₂. For some models Zero-emmission clusters removed during preprocessing and imputed for modeling consistency.

Another Visual Comparing CO2 and Energy Usage kWh over time.



Violin Plot

Illustrates the Evolving Distribution of CO₂ across Energy Usage Levels



- Initial Variable Distributions">Histograms revealed skewed CO₂ emissions and wide energy usage range. Guided transformation choices (e.g., log_kWh) and flagged missing values for imputation.

The violin plot reveals how CO₂ emission distributions change at different energy consumption levels, showing both frequency and probability density.

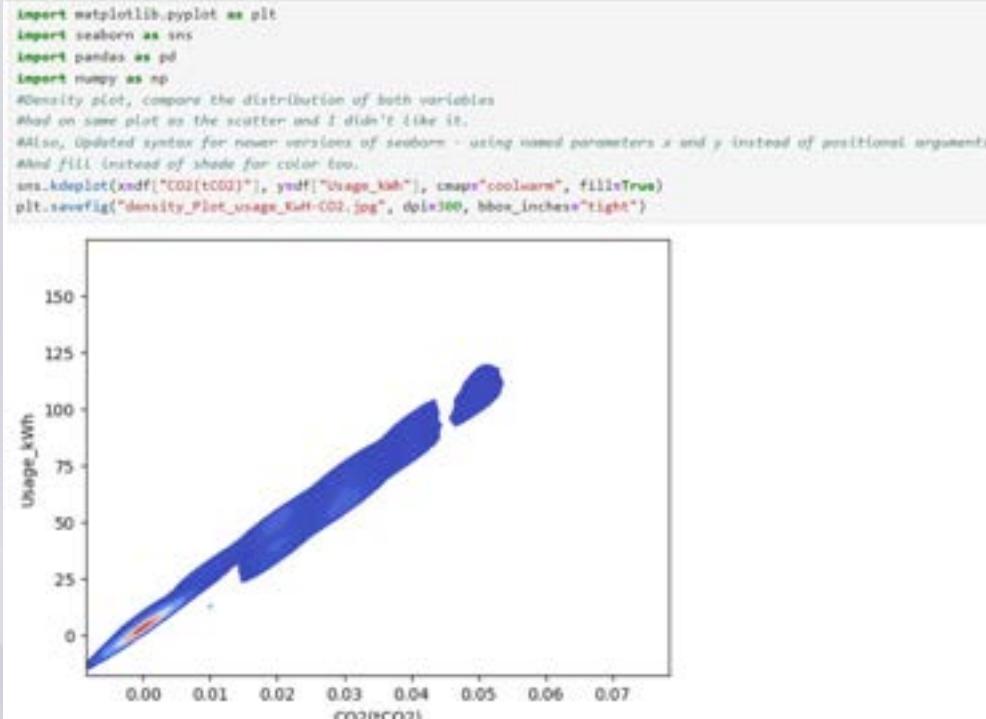
The plot revealed a dense spike at zero CO₂ emissions with a surprising vertical tail stretching into high energy usage levels. These outliers indicated likely data entry gaps since high kWh should yield measurable emissions. This informed targeted imputation to restore data integrity ahead of modeling

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

```
load('train.csv')
# Load libraries
library(tidyverse)
library(ggplot2)
library(violinette)

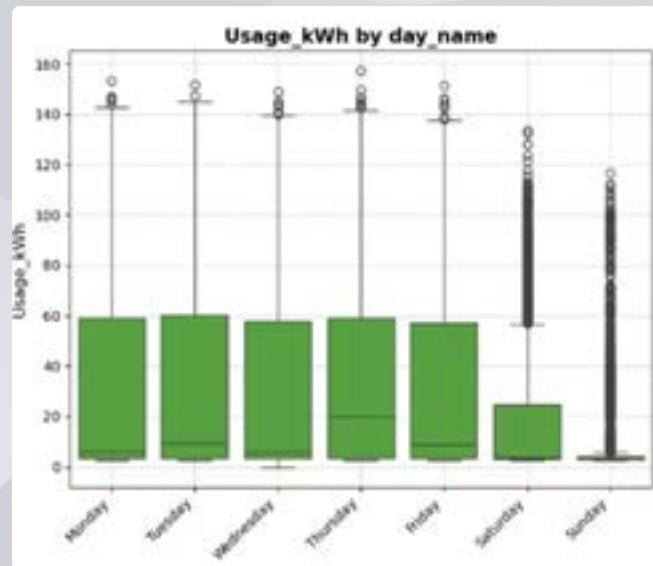
# Load the violin
violin::ViolinPlot(usage_kWh ~ CO2/CO2, data = train,
                   pch.outliers = 15, outlier.numb = 10, outlier.size = 10, outlier.colour = "#00FFFF",
                   pch.outline = 15, outline.numb = 10, outline.size = 1, outline.colour = "#000000",
                   pch.mean = 15, mean.size = 1, mean.colour = "#000000",
                   pch.violin = 15, violin.size = 1, violin.colour = "#000000",
                   pch.outline = 15, outline.numb = 10, outline.size = 1, outline.colour = "#000000")
```

Looking at the Density (Kernel)Plot to look



- The predominantly blue color with a small amount of red around intercept CO2, suggests that the highest density (red) is concentrated near the lower end of the CO2 axis, with density decreasing (shifting to blue) as CO2 increases.
- The thicker density from 0.02 to 0.05 CO2, 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 CO2 (x-axis) and energy usage (y-axis).

Boxplots illustrating how Energy Usage Fluctuates across time



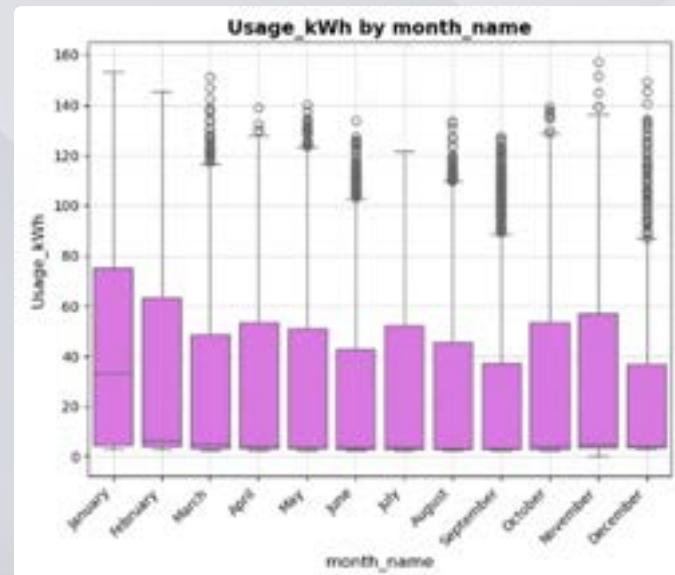
Day of the Week

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.



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

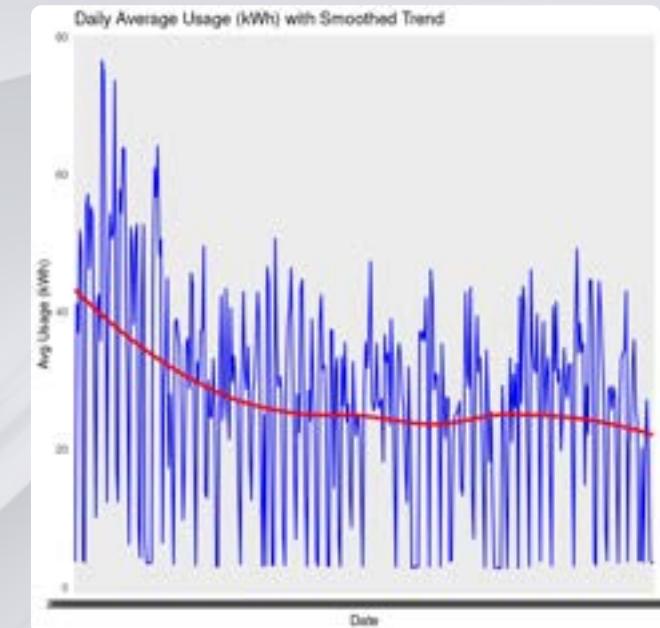
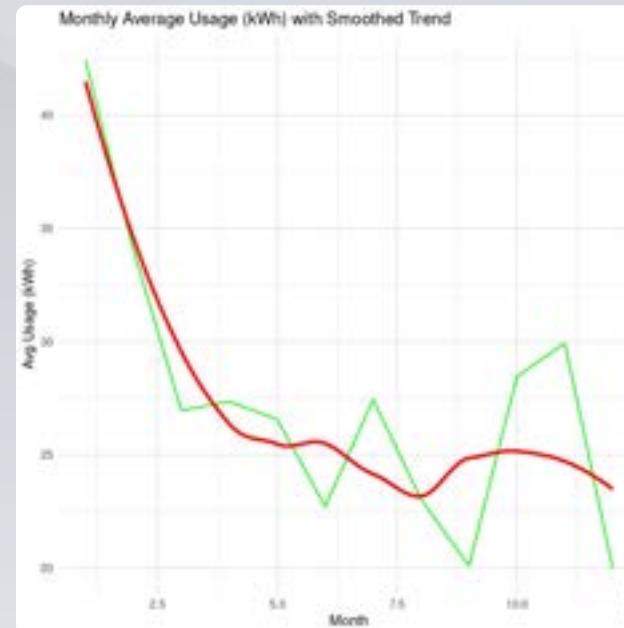
Supplements to the Boxplots showing monthly trends

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.

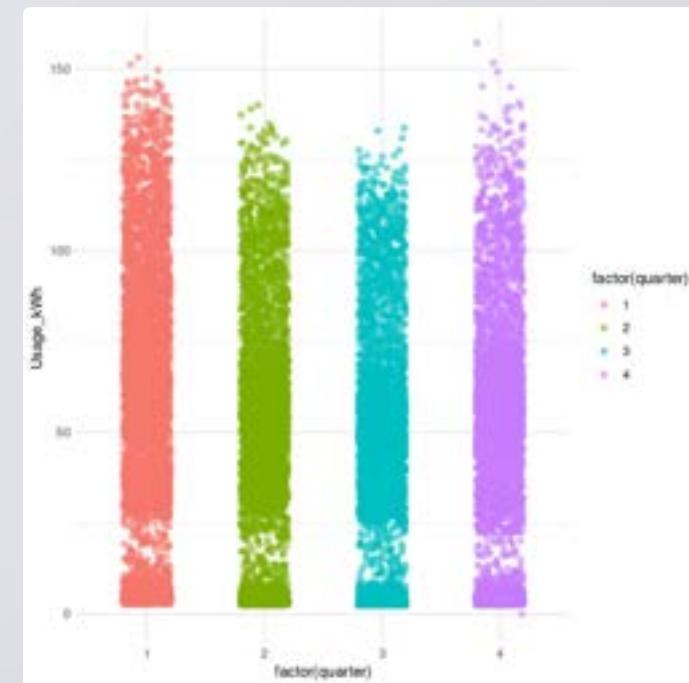
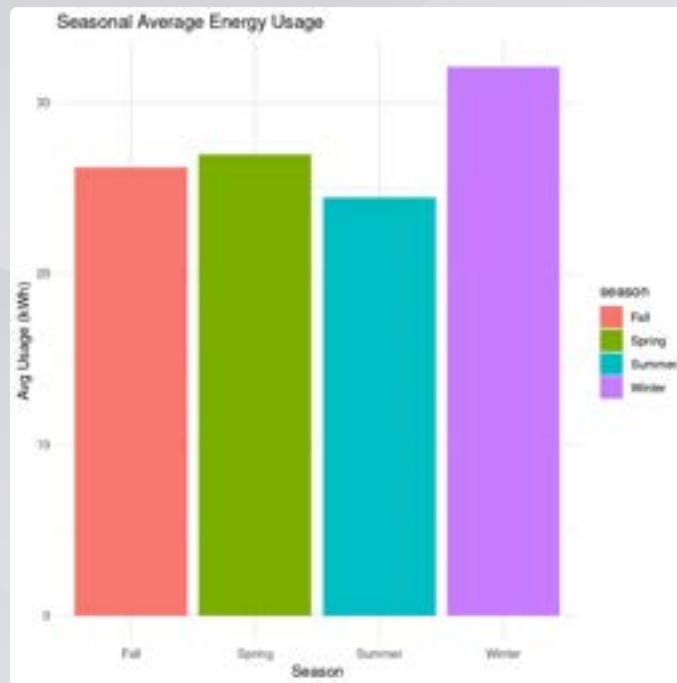
CODE USING R IN KAGGLE

```
Plotting of the monthly trend as well as time lagged average
```

```
average_usage = lagged_avg_usage %>%
```



SEASONS and Energy Usage (kWh)



```
# Load all seasonal intervals, map quarter and month to them, but just for big_observ
# quarters
# Note: c(1, 2, 3) = Fall, Spring, Summer
# Note: c(4, 5, 6) = Winter
# Note: c(7, 8, 9) = October, November, December
# Note: c(10, 11, 12) = January, February, March
```

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

```
#geom_jitter(. , size = factor(1), alpha = 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	5.489759e+07	1.0	2.145139e+07	1.642754e-06
relative_humidity_2m	2.173542e+07	1.0	5.399054e+07	7.346798e-07
precipitation	2.787166e+07	1.0	2.447770e+07	2.338911e-07
CO2	2.338833e+07	1.0	2.340871e+07	8.888888e-08
Residual	9.140417e+07	2945.0	Null	Null

ANOVA For Energy Usage kWh without CO2

	sum_sq	df	F	PR(>F)
temperature_2m	6.205386e+07	1.0	608.144358	3.787712e-121
relative_humidity_2m	2.001124e+07	1.0	2913.916111	0.000000e+00
precipitation	4.107116e+07	1.0	48.017838	1.487520e-10
Residual	5.459986e+07	2945.0	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 CO₂ is included:

- CO₂ 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 CO₂.

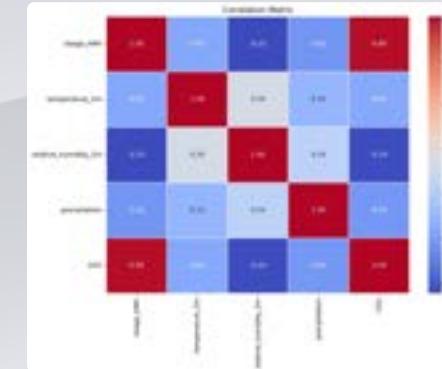
When CO₂ is excluded:

- Relative humidity** becomes the strongest predictor ($F= 2914$).
- Temperature** also has a sizable effect ($F= 608$).
- 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:

CO₂ 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 CO₂ with Weather-related environmental variables — including temperature, humidity, and precipitation.

While the correlation between CO₂ 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, datadf).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, datadf).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	8.159554	1.0	662.281719	1.058104e-144
relative_humidity_2m	0.675527	1.0	2884.004621	0.000000e+00
precipitation	0.009517	1.0	39.503582	3.313068e-10
Residual	8.440702	35036.0	NaN	NaN

CO₂ as a Response Variable

- **Usage_kWh is once again the strongest driver of CO₂ 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

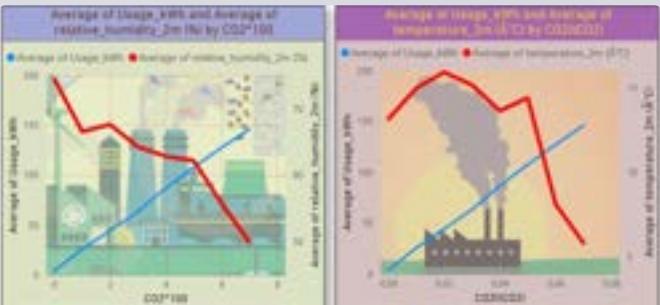
ANOVA CHECK

```
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 model
model = smf.ols(formula, dataset).fit()
anova_table = sm.stats.anova_lm(model, typ=2)

print(anova_table)
```



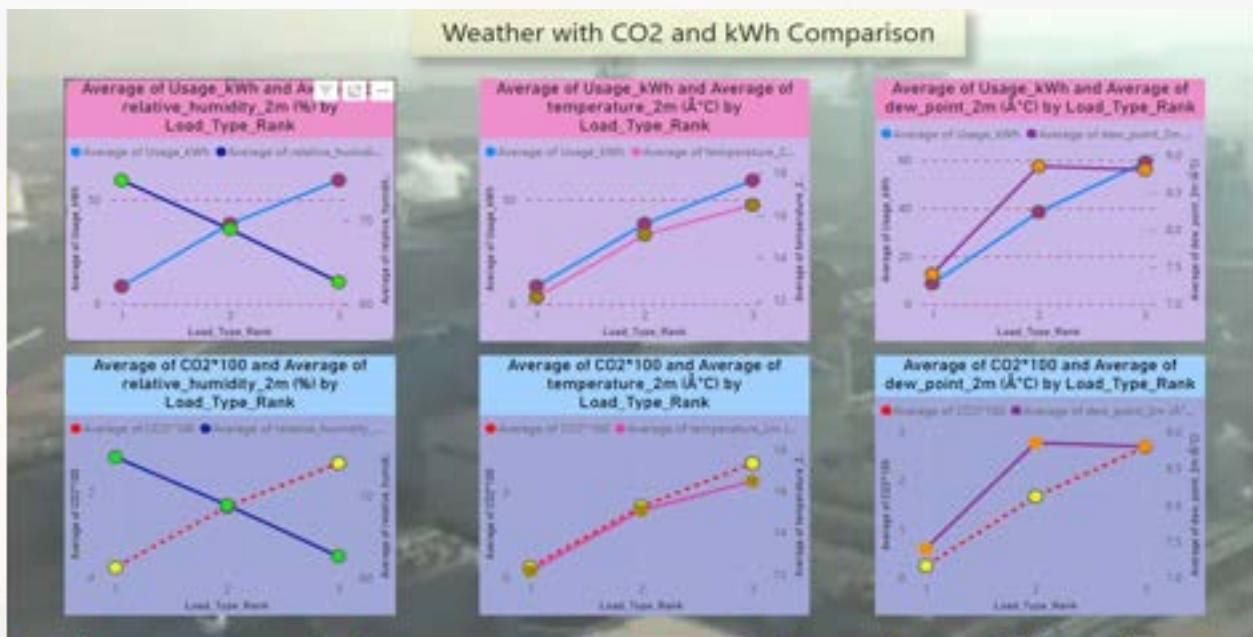
	sum_sq	df	F \
Usage_kWh	8.226553e+00	1.0	1.348491e+06
temperature_2m	4.476525e-04	1.0	7.337892e+01
Usage_kWh:temperature_2m	3.474503e-05	1.0	5.695383e+00
relative_humidity_2m	2.218008e-05	1.0	3.635745e+00
Usage_kWh:relative_humidity_2m	4.227809e-04	1.0	6.930199e+01
precipitation	2.653875e-07	1.0	4.350216e-02
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 CO₂ variation
- Interactions with temperature and humidity are crucial:** they reveal that energy usage affects CO₂ differently depending on these environmental conditions.
- Precipitation seems negligible here**, both on its own and in interaction with Usage_kWh.

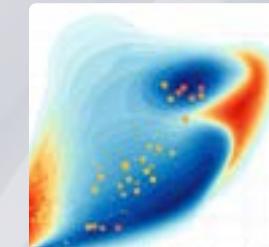
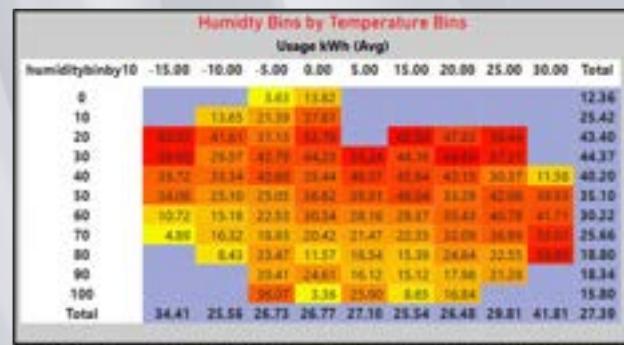
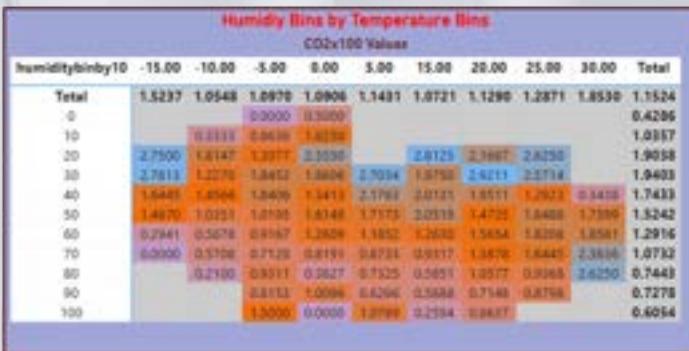
Section 3: Feature Engineering and Diagnostics



Key Points:

- Energy and CO₂ 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

HEATMAP of the HumidityxTemp Interaction for both CO2 and Usage kWh



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

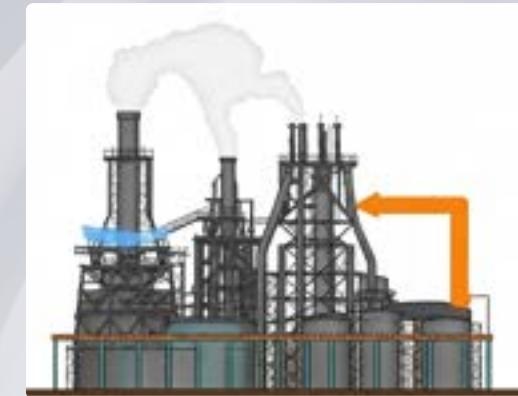
Measures the extent to which current leads voltage in a capacitive circuit.

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)

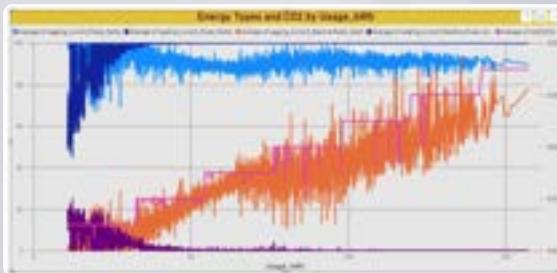
Represents the total amount of **real energy consumed** by the steel plant over time while capturing the core operational demand across all machinery and processes.

Load Rank (1–3)

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

Energy Variables Graphs

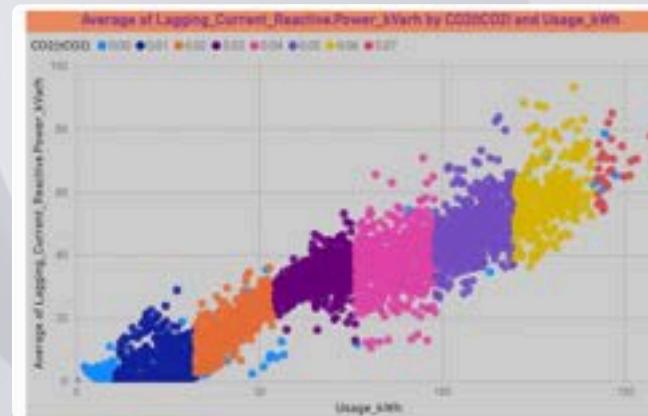
Multi-Variable Energy Scatter: Mapping Power Behavior and CO₂ Emissions



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

The CO₂ follows the same trend and is on top of the Lagging Reactive Power.

Energy Usage, Reactive Load, and CO₂



Color gradations represent CO₂ intensity, with warmer hues signaling rising emissions

The **stacked segments** hint at discrete operational modes and each with its own CO₂ footprint

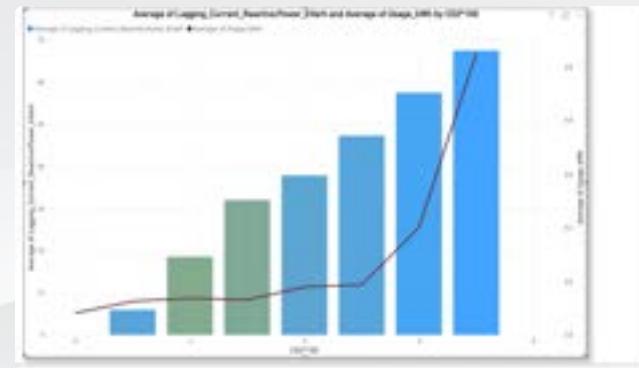
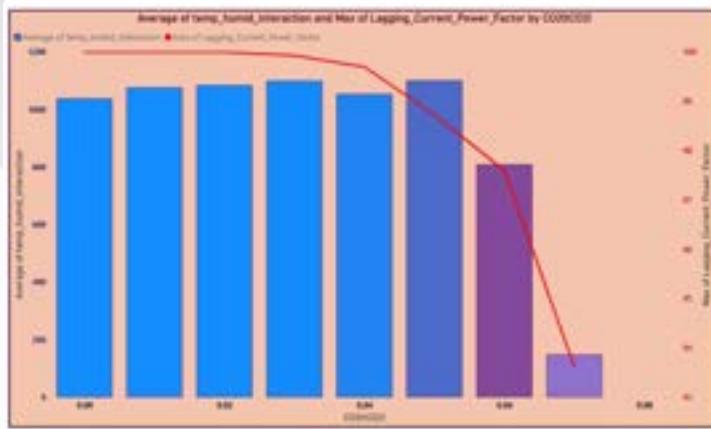
Outliers in light blue scattered throughout may be imputed gaps and unexplained behavior.

The Visual of Steelplant Dynamics

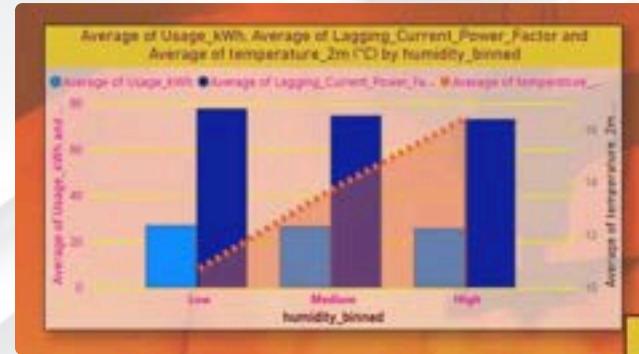
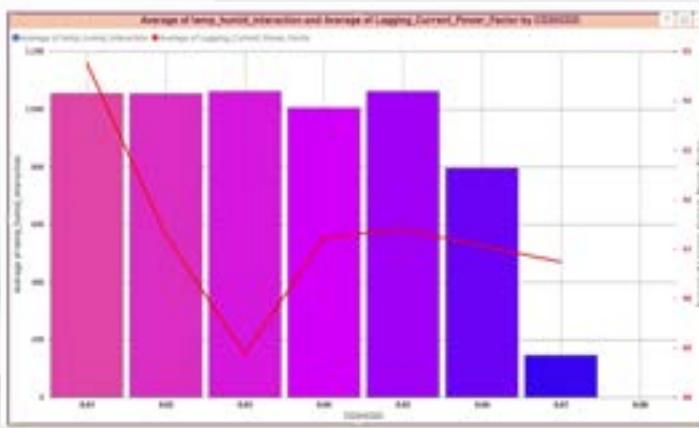
Each color stratum performs its part in the energy-emissions ensemble. The two Scatter plots charts **Usage kWh** against **Lagging Reactive Power kVarh**, forming a linear composition. The **color-coded bands**, segmented by CO₂ levels from light blue (0.0) to red (~0.07), reveal distinct operational layers across the plant's load profile.

Lagging Current Reactive Power kVarh variable

Energy Behaviors Leads Emission Shifts with environmental weather influence

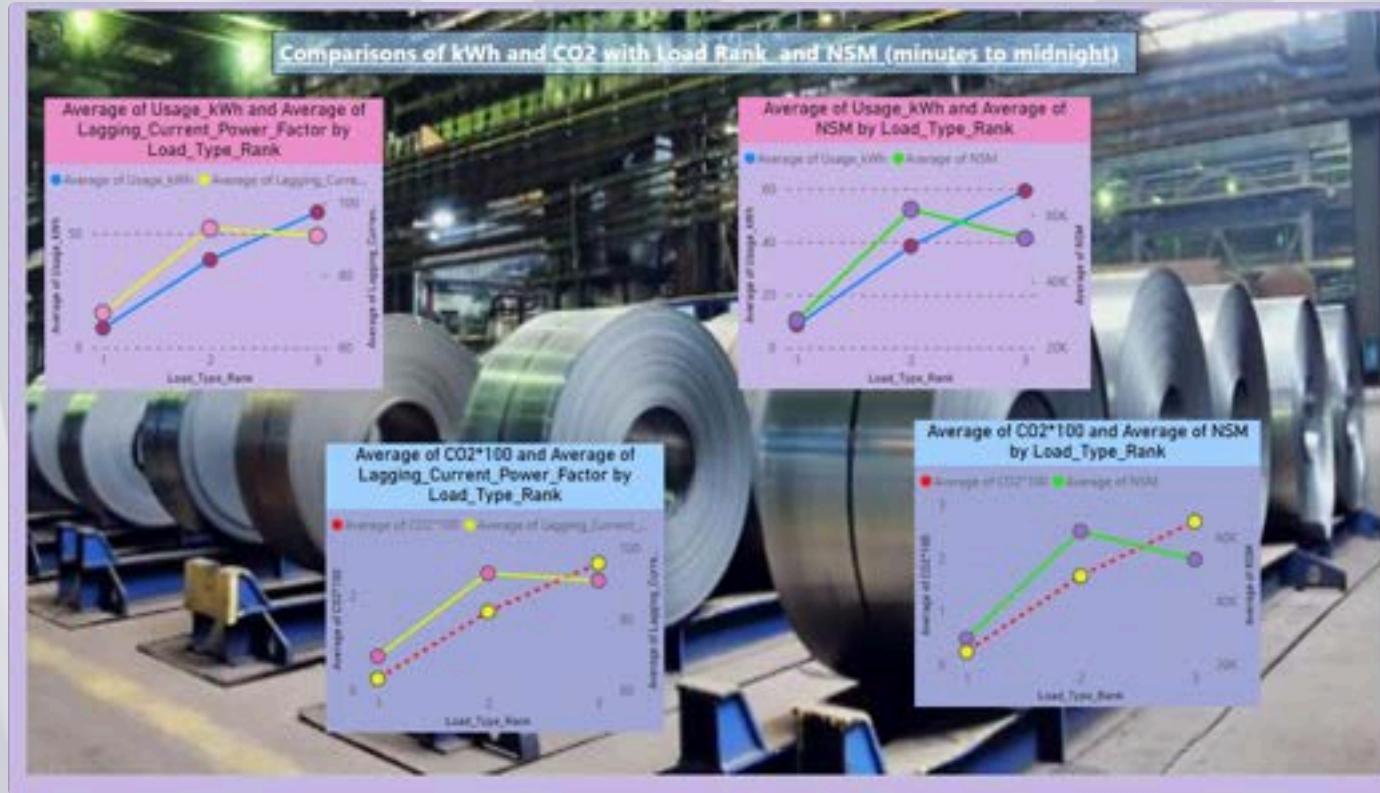


Both lagging energy and usage rise with CO₂, but there's a *nonlinear surge* at the highest emission levels.



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 CO₂, 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 CO₂
- High CO₂ 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)

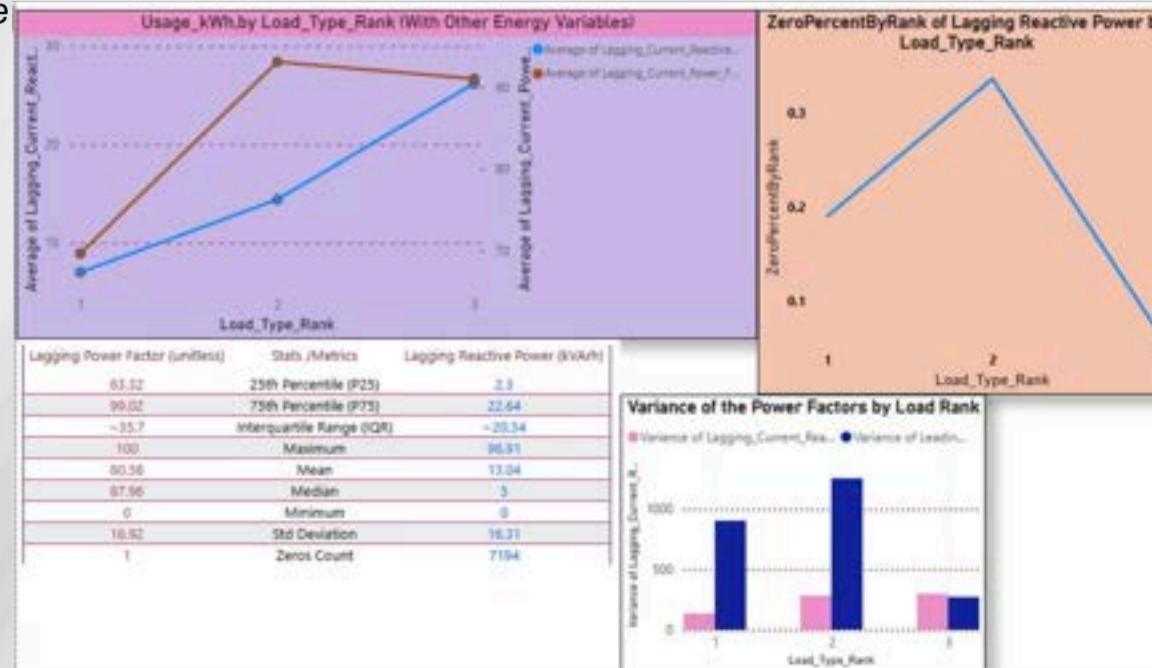
Definition: Total reactive energy used (inductive demand)

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 CO₂

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

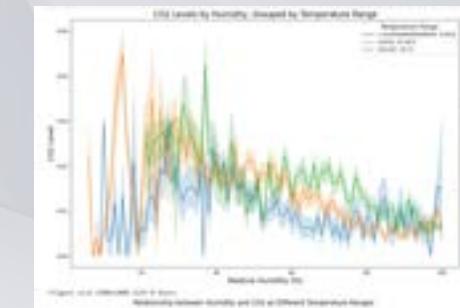
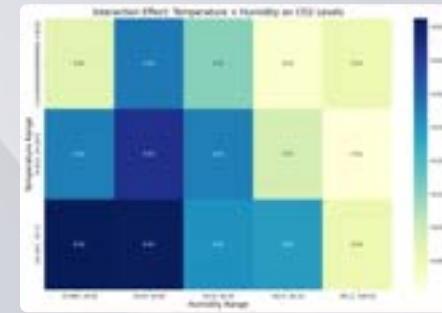


Interpretation Notes

- Model explains ~78.9% of variance in CO₂
- The interaction effect is powerfully negative, indicating that high humidity dampens the temperature-driven CO₂ 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 CO₂ 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 ₂ offset
Temp (°C)	3.119e-05	< 0.001	↑ Temp → slight increase in CO ₂
Reactive Power (kVarh)	0.0009	≈ 0	Strong positive CO ₂ association
Temp × Humidity Interaction	-4.63e-06	≈ 0	Highly significant suppressive effect



OLS Regression Model Continues for CO₂

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

Will multicollinearity be reduced with added the different Energy variables?

```
# Step 2: Add power variables and see if multicollinearity is reduced
power_vars = c("Leading_Current_Reactive.Power_kWark",
              "Leading_Current_Power_Factor",
              "Lagging_Current_Reactive.Power_kWark",
              "Lagging_Current_Power_Factor",
              "Leading_Current_Power_Factor")
```

Final Model Summary:

OLS Regression Results						
Dep. Variable:	CO2(kgCO2)	R-squared:	0.897			
Model:	OLS	Adj. R-squared:	0.897			
Method:	Least Squares	F-statistic:	5.106e+04			
Date:	Fri, 28 Jun 2023	Prob (F-statistic):	0.00			
Time:	18:27:53	Log-Likelihood:	-1.3474e+05			
Nb. Observations:	25648	AIC:	-2.695e+05			
Df Residuals:	25643	BIC:	-2.684e+05			
Df Model:	6					
Covariance Type:	opg					
	coef	std err	t	P> t	[9.025	0.975]
const	-0.001	0.000	-1.07.096	0.000	-0.000	-0.000
temperature_2m (°C)	-0.0001	1.46e-05	-9.898	0.000	-0.000	-0.45e-05
relative_humidity_2m (%)	-5.421e-05	2.48e-06	-22.247	0.000	-5.7e-05	-0.95e-05
temp_humid_interaction	2.488e-05	1.55e-07	16.295	0.000	2.15e-05	2.78e-05
Lagging_Current_Reactive.Power_kWark	0.0007	2.12e-06	322.749	0.000	0.001	0.001
Lagging_Current_Power_Factor	0.0004	2.05e-06	179.383	0.000	0.000	0.000
Leading_Current_Power_Factor	0.0002	1.32e-06	123.627	0.000	0.000	0.000
Omnibus:	4713.063	Durbin-Watson:	0.529			
Prob(Omnibus):	0.000	Sargan-Bera (B):	83639.818			
Skew:	0.183	Prob(JB):	0.88			
Kurtosis:	9.494	Cond. No:	1.44e+04			

Notes:

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

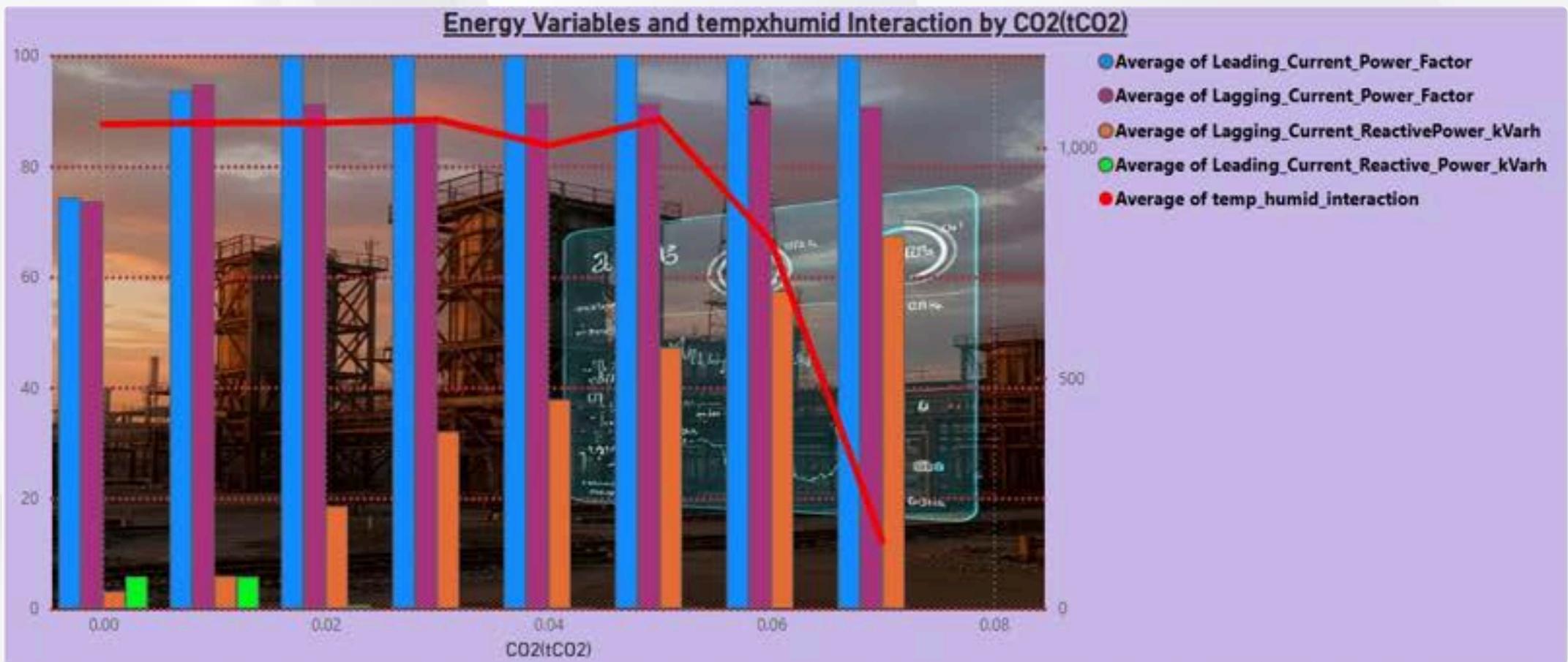
[2] The condition number is large, 1.44e+04. This might indicate that there are strong multicollinearity or other numerical problems.

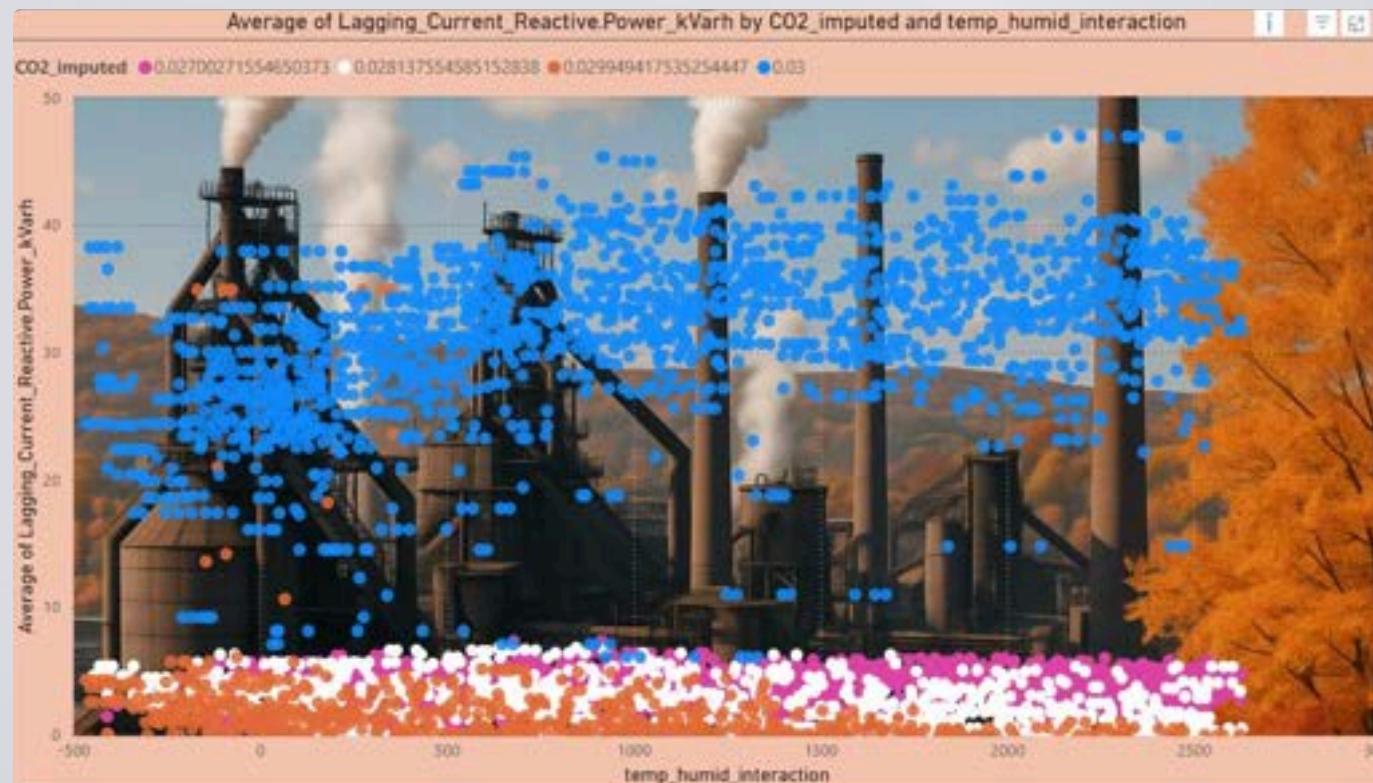
Variance Inflation Factors for Final Model:		
Variable	VIF	
const	117.730342	
temperature_2m (°C)	14.726059	
relative_humidity_2m (%)	3.059169	
temp_humid_interaction	19.660000	
Lagging_Current_Reactive.Power_kWark	1.577887	
Lagging_Current_Power_Factor	1.959991	
Leading_Current_Power_Factor	2.312516	

Interpretation

- The model explains ~90% of the variance in CO₂, 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.

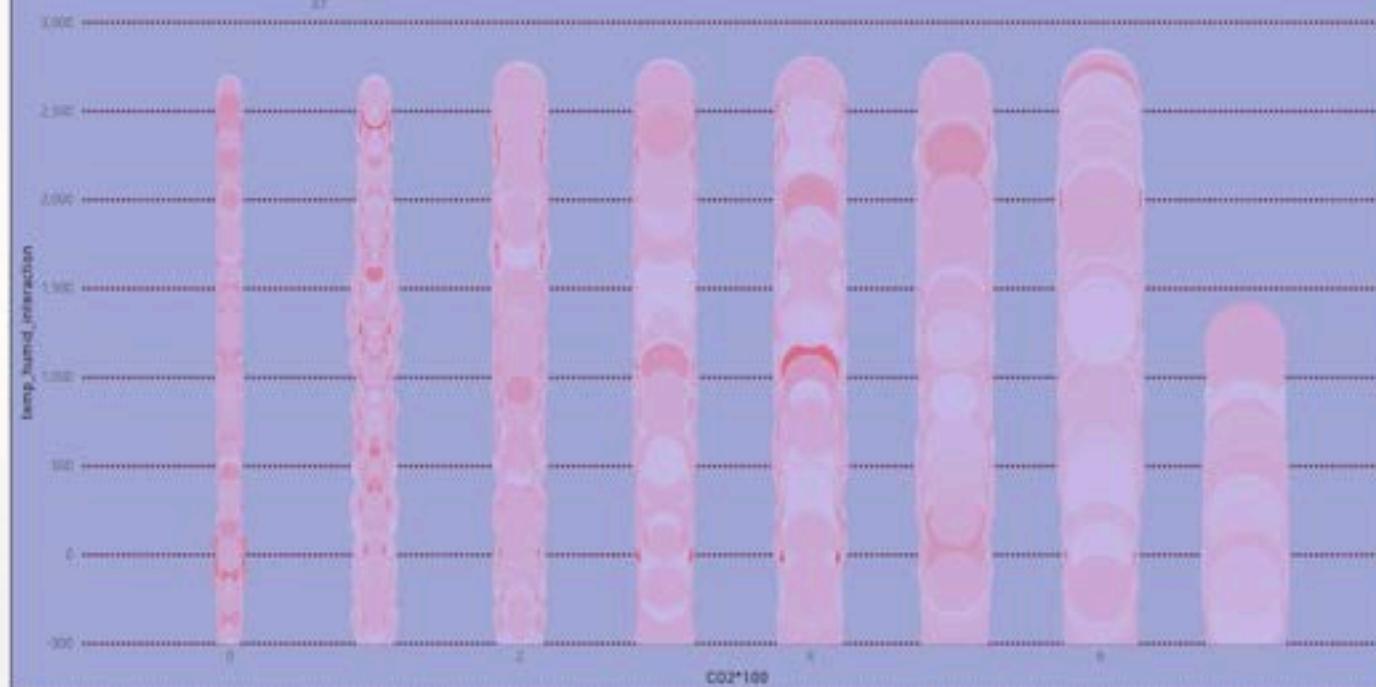
Graphing of the variables in the CO2 Model





Average of Lagging_Current_ReactivePower_kVArh by CO2*100 and temp_humid_interaction

Count of temp_humid_interaction: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71



OLS Regression Model Continues for Energy Usage kWh Following the same model as previous CO₂

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:	998.1
Date:	Fri, 28 Jun 2019
P-value (F-statistic):	0.00
Time:	13:27:53
Log-Likelihood:	-1.712e+05
No. Observations:	35648
AIC:	3.425e+05
DF Residuals:	35638
BIC:	3.426e+05
DF Model:	3
Covariance Type:	opg
coef std err t P> t [0.025 0.975]	
const 53.9243 0.341 15.565 0.000 50.279 52.769	
temperature_3m (°C) -0.8123 0.005 -12.442 0.000 -0.584 0.319	
relative_humidity_3m (%) -0.4259 0.015 -28.848 0.000 -0.455 -0.397	
temp_humid_interaction -0.8951 0.001 -5.783 0.000 -0.997 -0.994	
Observations: 5439.568 Durbin-Watson: 0.191	
Prob> Omnibus : 0.000 Durbin-Watson (DW): 0.019468	
Skew: 1.181 Prob>F: 0.00	
Kurtosis: 5.631 Cond. No.: 7.28e+05	

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 7.3e+05. This might indicate that there are strong multicollinearity or other numerical problems.

Results of adding each power variable with p <= 0.10:

Variable	AIC	R-squared
0 Lagging_Current_Reactive_Power_WkWh	281585.626460	0.019718
1 Lagging_Current_Reactive_Power_WkWh	285362.129874	0.029967
2 Lagging_Current_Power_Factor	285186.516853	0.000915
3 Lagging_Current_Power_Factor	288115.346893	0.012375

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:	5.213e+04			
Date:	Fri, 28 Jun 2019	Prob > F-statistic:	0.00			
Time:	13:14:07	Log-Likelihood:	-1.388e+05			
No. Observations:	35648	AIC:	2.685e+05			
DF Residuals:	35642	BIC:	2.682e+05			
DF Model:	3					
Covariance Type:	opg					
coef std err t P> t [0.025 0.975]						
const -88.3542 0.737 -59.369 0.000 -81.998 -79.111						
temperature_3m (°C) -0.2154 0.008 -26.618 0.000 -0.255 -0.176						
relative_humidity_3m (%) -0.1099 0.003 -23.576 0.000 -0.119 -0.101						
temp_humid_interaction 0.0048 0.000 16.217 0.000 0.004 0.005						
lagging_Current_Reactive_Power_WkWh 1.4679 0.004 357.526 0.000 1.466 1.470						
lagging_Current_Reactive_Power_WkWh 0.3667 0.003 10.518 0.000 0.323 0.418						
lagging_Current_Power_Factor 0.7395 0.004 186.887 0.000 0.732 0.747						
lagging_Current_Power_Factor 0.3079 0.005 71.246 0.000 0.307 0.400						
Observations: 35648 Durbin-Watson: 0.378						
Prob> Omnibus : 0.000 Durbin-Watson (DW): 0.0007169						
Skew: 0.784 Prob>F(3,34): 0.00						
Kurtosis: 6.143 Cond. No.: 1.35e+04						

Variance Inflation Factors for Final Model:	
Variable	VIF
const	103.914989
temperature_3m (°C)	14.795932
relative_humidity_3m (%)	0.399613
temp_humid_interaction	29.660012
lagging_Current_Reactive_Power_WkWh	1.681398
lagging_Current_Reactive_Power_WkWh	0.657828
lagging_Current_Power_Factor	0.895183
lagging_Current_Power_Factor	0.705265



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

OLS Regression for Energy Usage kWh

Model Performance Interpretation



Environment Variables Only

R² Score: 0.078

AIC: 342,600

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



Add Power Factors Incrementally

R² Score: Up to 0.912

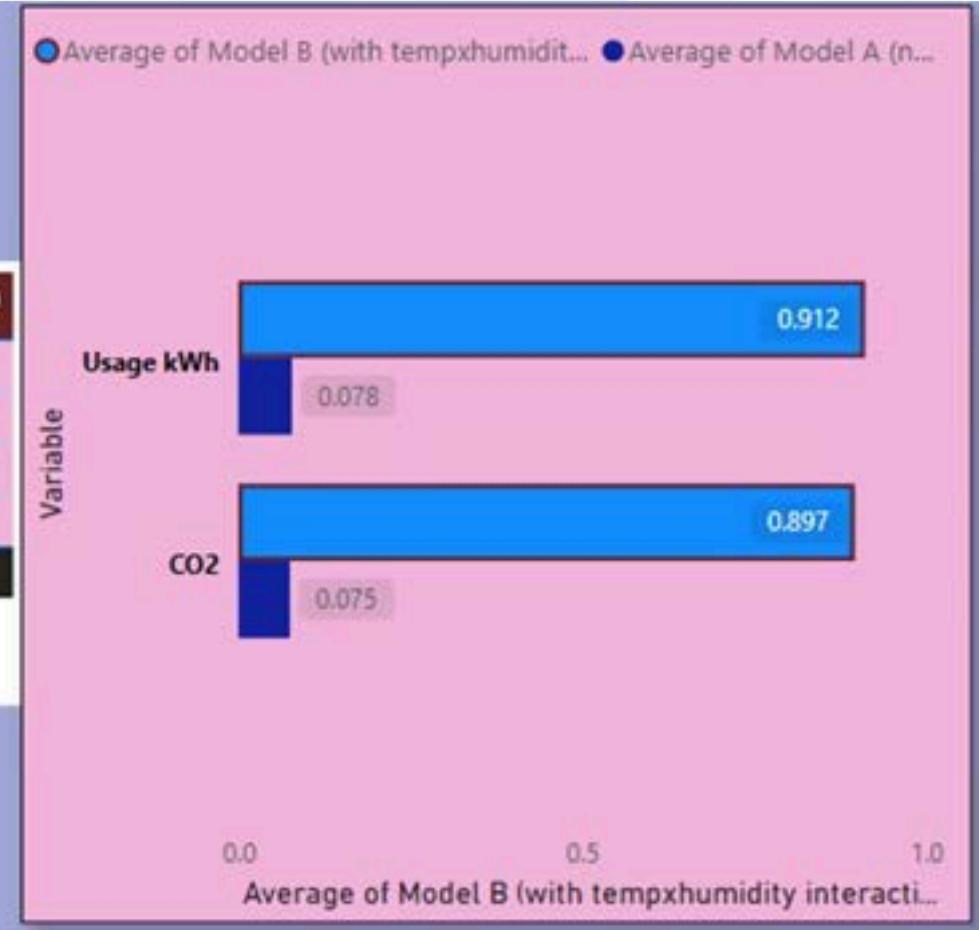
AIC: ↓ to 260,100

Dramatic boost — power dynamics dominate energy usage patterns

Adding power variables to the environmental baseline yields a 12x increase in explanatory strength (R² 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 ² (CO ₂)	0.075	0.897
R ² (Usage_kWh)	0.078	0.912
AIC (CO ₂)	-192,439	-269,194
AIC (Usage_kWh)	342,587	260,115
Total		



OLS Regression Model for CO₂

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(CO2)	R-squared:	0.879			
Model:	OLS	Adj. R-squared:	0.879			
Method:	Least Squares	F-statistic:	3499.			
Date:	Mon, 30 Jun 2025	Prob (F-statistic):	0.00			
Time:	11:18:07	Log-Likelihood:	-9812.1			
No. Observations:	35040	AIC:	-19624.2			
DF Residuals:	35034	BIC:	-19542.1			
DF Model:	5					
Covariance Type:	opg					
<hr/>						
	coef	std err	t	P> t	[95.0%	8.379]
Intercept	8.0000	0.0000	81.000	0.0000	8.004	8.023
temperature_2m (°C)	8.0000	0.01e-06	21.818	0.0000	8.000	8.000
relative_humidity_m (%)	-8.0000	4.43e-06	-1.812	0.0000	-8.000	-8.000
<hr/>						
Residuals:	6044.560	Satterthwaite:	0.207			
F-statistic:	8.0000	Auger-Bera (W):	7081.007			
Df Residuals:	3499	Df Model:	5			
Df Model:	2,440	Cond. No.:	276.			
<hr/>						
Refined:						
(1) Standard Errors assume that the covariance matrix of the errors is correctly specified.						
R-squared: 0.8794172675548825						
AIC: -19624.4508153099						
Power variables found in dataframe: ["Lagging_Current_Reactive_Power_Mvarh", "Leading_Current_Reactive_Power_N Varh", "Lagging_Current_Power_Factor", "Leading_Current_Power_Factor"]						

Results of adding each power variable:						
Variable	AIC	R-squared	%	Coefficient	P-value	
Lagging_Current_Reactive_Power_Mvarh	-246640.152993	0.882071				
Leading_Current_Reactive_Power_Mvarh	-246640.611165	0.882073				
Lagging_Current_Power_Factor	-256644.948135	0.812894				
Leading_Current_Power_Factor	-269194.218847	0.896105				

Final Model Summary:						
OLS Regression Results						
Dep. Variable:	CO2(CO2)	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:	11:18:07	Log-Likelihood:	1.3466e+05			
No. Observations:	35040	AIC:	-2.692e+05			
DF Residuals:	35034	BIC:	-2.691e+05			
DF Model:	5					
Covariance Type:	opg					

Note: (1) Standard Errors assume that the covariance matrix of the errors is correctly specified.
(2) The condition number is large, 1.0000. This might indicate that there are strong multicollinearity or other numerical problems.
Model select: AIC=692.6000000000000, BIC=696.6000000000000

Influence Statistics for Final Model:

variable	influence	residuals
Intercept	0.0000	0.0000
temperature_2m (°C)	0.0000	0.0000
relative_humidity_m (%)	0.0000	0.0000
Lagging_Current_Reactive_Power_Mvarh	0.0000	0.0000
Leading_Current_Reactive_Power_Mvarh	0.0000	0.0000
Lagging_Current_Power_Factor	0.0000	0.0000
Leading_Current_Power_Factor	0.0000	0.0000



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

FINAL OLS Regression Model for CO₂

Without tempxhumid interaction and Leading Current Reactive Power

Final Model Summary: OLS Regression Results						
Dep. Variable:	CO2(kg/kWh)	R-squared:	0.897			
Model:	OLS	Adj. R-squared:	0.897			
Date:	Fri, 30 Jun 2023	F-statistic:	6.87e+04			
Time:	13:14:00	Prob (F-statistic):	0.00			
N Observations:	19848	(log-)Likelihood:	1.3408e+05			
Df Residuals:	19834	AIC:	-2.682e+05			
Df Model:	5	BIC:	-2.601e+05			
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.0397	0.000	-138.700	0.000	-0.048	-0.030
temperature_3m (°C)	0.0034e-05	3.84e-06	39.880	0.000	5.40e-05	6.58e-05
relative_humidity_3m (%)	-2.463e-05	1.63e-06	-15.895	0.000	-2.79e-05	-2.14e-05
Lagging_Current_Reactive_Power_kVarh	0.0007	2.14e-06	321.354	0.000	0.002	0.001
Lagging_Current_Power_Factor	0.0004	2.05e-06	178.294	0.000	0.000	0.000
Leading_Current_Power_Factor	0.0002	1.12e-06	122.801	0.000	0.000	0.000
Observations:	4729.682	Omnibus (Warning):	0.515			
Prob(Omnibus):	0.600	T-sqrd (Bera):	0.027.418			
Skew:	0.117	Prob(z):	0.00			
Kurtosis:	9.418	Cond. No.:	1.48e+01			
Notes:						
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.						
[2] The condition number is large, 1.63e+03. This might indicate that there are strong multicollinearity or other numerical problems.						
Final model R-squared: 0.8965953598042129						
Final model AIC: -268194.21834737787						
Variance Inflation Factors for Final Model:						
	Variable	VIF				
0	const	186.810198				
1	temperature_3m (°C)	1.207647				
2	relative_humidity_3m (%)	1.337503				
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 ²	AIC	Key Insight
Environment-only (No Interaction)	0.075	-192,440	Low explanatory power — temp & humidity contribute separately but weakly
Final CO ₂ 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

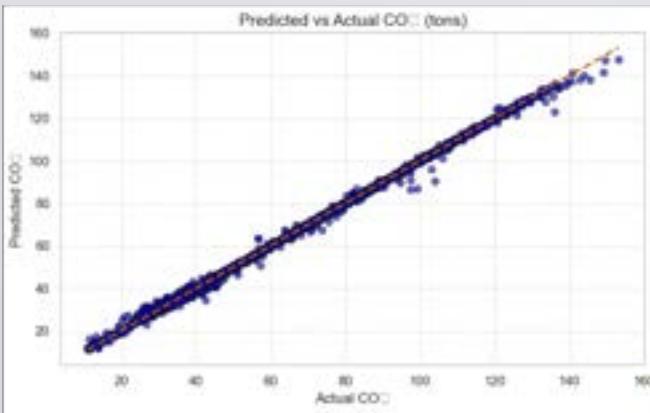
Interpretation

- 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 CO₂ 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

```
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()
```



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 CO₂ 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	---	---	---	---	---

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	---
I	-0.475	---	---	---	---	---
Omnibus:	5405.050	Durbin-Watson:	0.193	---	---	---
Prob(Omnibus):	0.000	Jarque-Bera (JB):	8245.949	---	---	---
Skew:	1.147	Prob(JB):	0.00	---	---	---
Kurtosis:	3.620	Cond. No.:	276.	---	---	---
Notes:						
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.						
R-squared: 0.07756887262074719						
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']						
Results of adding each power variable:						
	Variable	AIC	R-squared	\\		
0	Lagging_Current_Reactive_Power_kVarh	285387.734992	0.039715			
1	Leading_Current_Reactive_Power_kVarh	285364.892743	0.019843			
2	Lagging_Current_Power_Factor	265316.215698	0.098343			
3	Leading_Current_Power_Factor	260378.472452	0.911709			
Coefficients:						
0	Coefficient	P-value				
0	1.797456	0.000000e+00				
1	0.057878	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:	6.029e+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	---	---	---	---	---

S	0.975]	coef	std err	t	P> t	[0.02
const	83.5150	0.717	-116.558	0.00	---	---
temperature_2m (°C)	0.1813	0.006	27.307	0.00	---	---
relative_humidity_2m (%)	-0.0533	0.003	-17.056	0.00	---	---
I	-0.059	-0.047	---	---	---	---
Omnibus:	3944.361	Durbin-Watson:	0.279	---	---	---
Prob(Omnibus):	0.000	Jarque-Bera (JB):	38034.462	---	---	---
Skew:	0.203	Prob(JB):	0.00	---	---	---
Kurtosis:	1.189	Cond. No.:	1.88e+03	---	---	---
Notes:						
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.						
[2] The condition number is large, 6.00000. This might indicate that there are strong multicollinearity or other numerical problems.						
Final model R-squared: 0.9120000000000001						
Final model AIC: 260378.472452						
Variance Inflation Factors for Final Model:						
0	const	180.35094	---	---	---	---
1	temperature_2m (°C)	1.377593	---	---	---	---
2	relative_humidity_2m (%)	1.377593	---	---	---	---
3	Lagging_Current_Reactive_Power_kVarh	0.601561	---	---	---	---
4	Leading_Current_Reactive_Power_kVarh	0.601561	---	---	---	---
5	Lagging_Current_Power_Factor	1.995872	---	---	---	---
6	Leading_Current_Power_Factor	1.994439	---	---	---	---

Key Insights and Interpretation

Variable	R ² +	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)

• $R^2 = 0.912$ | 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 ($R^2 = 91.2\%$), 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, 28 Jun 2025	Fprob (F-statistic):	0.00			
Time:	13:14:07	Log-Likelihood:	-2,98124e+05			
N Observations:	359888	AIC:	2,086e+05			
Df Residuals:	359818	BIC:	2,087e+05			
Df Model:	5					
Covariance Type:	opnrobust					
Variance Inflation Factors For Final Model:						
	Variable	VIF				
	const	1.00				
	temperature_2m (°C)	2.101647				
	relative_humidity_2m (%)	2.377543				
	lagging_Current_Reactive_Power_kVarh	2.573868				
	lagging_Current_Factive_Factor	2.649934				
	Leading_Current_Factive_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)	X 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				
Observations:	3593,794	Durbin-Watson:	8.00			
Prob(Durbin-Watson):	0.989	Jarque-Bera (JB):	17321.589			
Skew:	0.775	Prob(JB):	0.00			
Kurtosis:	6.406	Cond. No.:	3.43e+03			

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	-1.514958	-0.386175			
1	-1.514958	-0.386175			
2	-1.514958	-0.386175			
3	-1.514958	-0.386175			
4	-1.514958	-0.386175			

	Lagging_Current_Reactive.Power_kVarh	Lagging_Current_Power_Factor
0	-0.430510	-0.386175
1	-0.521952	-0.375713
2	-0.588278	-0.344264
3	-0.581198	-0.348899
4	-0.123458	-0.219517

	Leading_Current_Power_Factor
0	0.351388
1	0.513368
2	0.513368
3	0.513368
4	0.513368

	Unscaled Variables	temperature_2m ("C)	relative_humidity_2m (%)	Lagging_Current_Reactive.Power_kVarh	Lagging_Current_Power_Factor	Leading_Current_Power_Factor
0	13	40.8	2.85	73.21	100.0	
1	13	40.8	4.46	96.77	100.0	
2	13	40.8	3.08	76.29	100.0	
3	13	40.8	3.56	68.09	100.0	
4	13	40.8	4.30	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 score
rmse_reduced = np.sqrt(mse_reduced)
print(f"kWh Model (without Leading_Current_Power_Factor) - CV MSE: {mse_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: 1.2466929
```

Cross-Validation Results Summary

Model	CV MSE	CV RMSE	Interpretation
CO ₂ Model	0.000031	0.0055	Incredibly tight error and model generalizes well on CO ₂ 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 CO₂ and Energy Usage kWh

Overview: Where energy demand reflects an intricate balance of power quality and environmental conditions, CO₂ 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

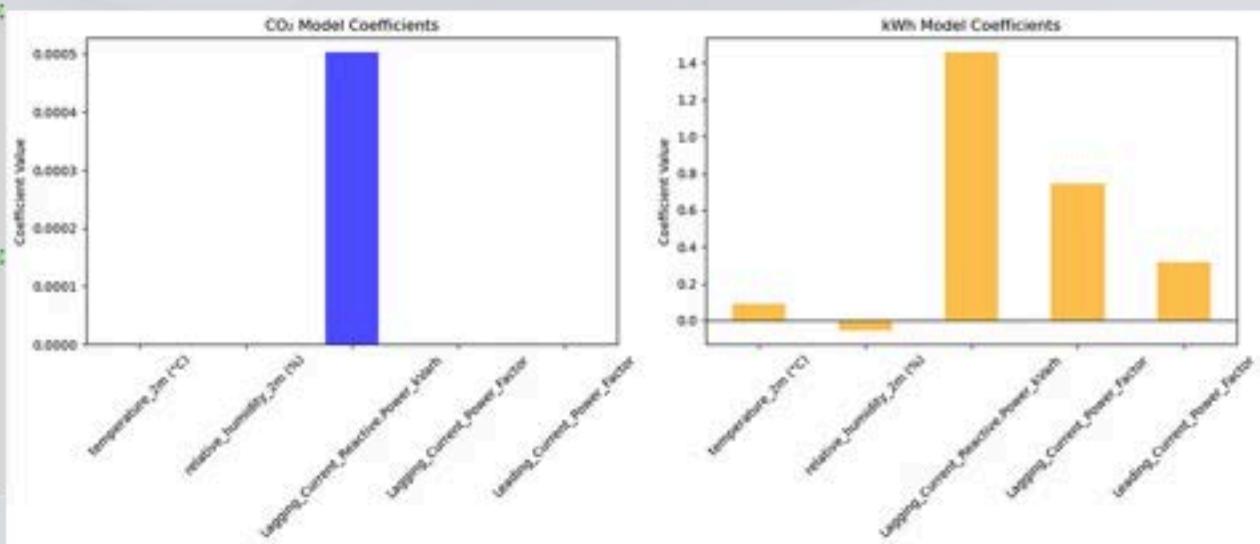
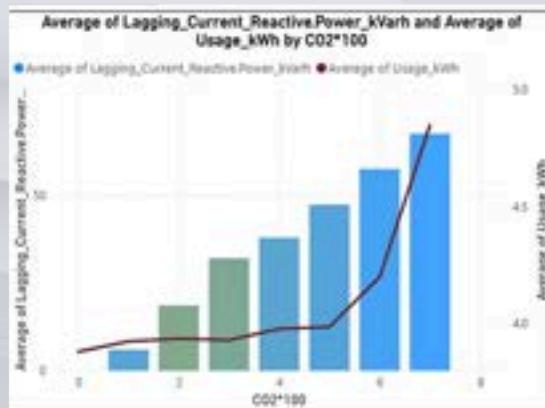
#Create the coefficient comparison dataframes
coeffs_df = pd.DataFrame()
'x_co2': lasso_co2.coef_,
'y_kwh': lasso_kwh.coef_,
), index=predictors)

#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 dataframes
coeffs_df = pd.DataFrame()
'x_co2': lasso_co2.coef_,
'y_kwh': lasso_kwh.coef_,
), index=predictors)
```

Previous Comparison of Lagging Reactive Power, CO2 and Energy Usage kWh



CO₂ Model

- Only Lagging Reactive Power survives LASSO penalty.
- That single non-zero bar is emissions proxy: if it draws VARh, it emits CO₂.
- 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]]

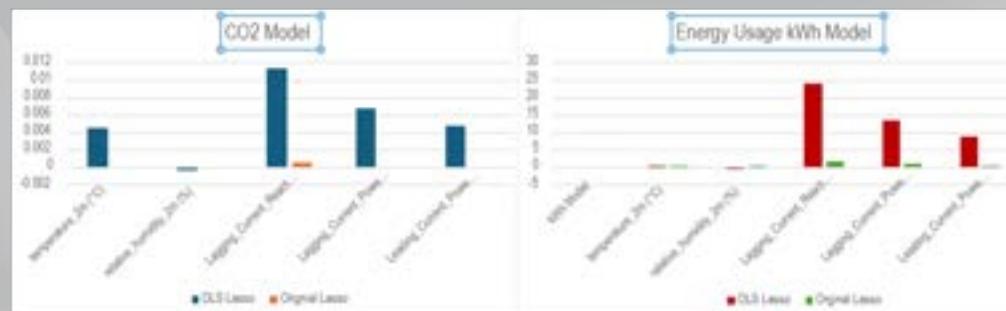
model_co2_lasso = sm.OLS(y_co2, lasso_selected_co2).fit()
model_kwh_lasso = sm.OLS(y_kwh, lasso_selected_kwh).fit()

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 (%)      -0.430668
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 R² (CO₂=0.897, kWh = 0.911) Cross-validation RMSEs remain tight (5-fold cross-validation)

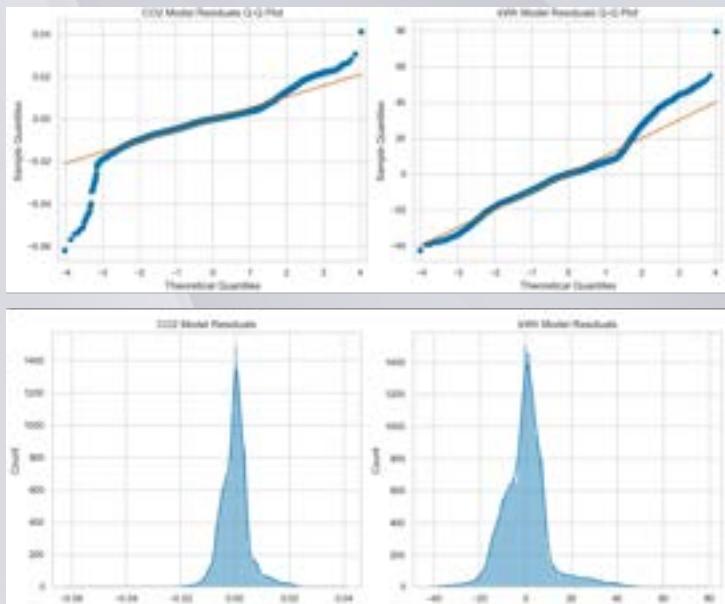
OLS Model (Lasso-selected variables)							OLS Model (Lasso-selected variables)								
OLS Regression Results							OLS Regression Results								
Dep. Variable:	CO2 (t/1000)	F-statistic:	6.03e+04	R-squared:	0.897	Adj. R-squared:	0.897	Dep. Variable:	Usage_kWh	F-statistic:	9.174e+04	R-squared:	0.911		
Model:	OLS	Prob (F-statistic):	0.000	Method:	Lasso	Adj. R-squared:	0.911	Model:	OLS	Prob (F-statistic):	0.000	Method:	Lasso	Adj. R-squared:	0.911
Date:	Mon, 20 Jun 2023	Log-Likelihood:	-5.3000e+05	Time:	11:16:18	AIC:	-1.6000e+05	Date:	Mon, 20 Jun 2023	Log-Likelihood:	-1.3012e+05	Time:	11:16:08	AIC:	-1.4012e+05
No. Observations:	75000	BIC:	-1.5820e+05	Df Residuals:	74994	BIC:	-1.4820e+05	No. Observations:	75000	BIC:	-1.4800e+05	Df Residuals:	74994	BIC:	-1.4810e+05
Df Model:	7			Df Model:	7			Df Model:	7			Covariance Type:	opg		
Convergence type:	converged			Covariance Type:	opg			Covariance Type:	opg						
	coef	std. err.	t	P> t	[0.025]	[0.975]		coef	std. err.	t	P> t	[0.025]	[0.975]		
const	9.4133	1.75e-05	5.31e-10	0.000	0.911	0.912	const	27.3869	0.853	32.349	0.000	27.282	27.491		
temperature_24h (°C)	0.0000	1.80e-05	11.000	0.000	0.801	0.802	temperature_24h (°C)	0.017	0.000	51.610	0.000	0.000	1.401		
relative_humidity_24h (%)	-0.0003	3.20e-05	-23.000	0.000	-0.801	-0.800	relative_humidity_24h (%)	-2.8768	0.000	-17.212	0.000	-1.399	-0.958		
lagging_current_reactive_Power_Mwh	9.9111	3.40e-05	2.81.354	0.000	0.911	0.913	lagging_current_reactive_Power_Mwh	21.7939	0.000	105.377	0.000	23.690	23.998		
lagging_Current_Power_Factor	0.0009	1.85e-05	1.01.200	0.000	0.807	0.807	lagging_Current_Power_Factor	14.8013	0.000	101.313	0.000	13.849	16.237		
leading_Current_Power_Factor	0.0009	1.83e-05	1.01.360	0.000	0.807	0.805	leading_Current_Power_Factor	9.4619	0.017	1.04.218	0.000	9.515	9.819		
Decimals:	4						Decimals:	3							
Prob(F statistic):	0.000						Prob(F statis):	0.000							
Stat:	0.517						Stat:	0.773							
Kurtosis:	5.814						Kurtosis:	6.876							

- It appears that CO₂ 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



CO₂ 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 CO₂, 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.

CO₂ 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 CO₂ 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 ₂ (tCO ₂)	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 ₂ (tCO ₂)	+1.15	Right-skewed → long tail of higher emissions
Usage_kWh	+1.20	Right-skewed → some extreme usage spikes

Both CO₂ 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 CO₂ 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 CO₂ 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 CO₂ 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 CO₂, due to its zero-inflation suspect, I couldn't resist testing what impact it might have on residual symmetry, skewness, and information criteria.

OLS Model (log-transformed): OLS Regression Results						
Deg. of Variables	log_CO2	R-squared	9.698			
Model:	0.6	Adj. R-squared:	9.698			
Method:	Least Squares	F-statistic:	6.208e-04			
Date:	Mon, 10 Jun 2023	P-value(F-statistic):	0.98			
Time:	15:14:12	Log-Likelihood:	3.3503e-05			
No. Observations:	39848	AIC:	-2.712e-05			
df Residuals:	39848	BIC:	-2.713e-05			
df Model:	9					
Covariance Type:	nonrobust					
coef std. err. t P> t [0.025 0.975]						
const	0.0113	2.7e-05	426.389	0.000	0.011	0.011
temperature_m (°C)	0.0006	2.96e-05	26.851	0.000	0.001	0.001
relative_humidity_m (%)	-0.0001	3.16e-05	-15.256	0.000	-0.001	-0.000
lagging_current_reactive_Power_kWh	0.0109	5.39e-05	123.442	0.000	0.011	0.011
lagging_current_power_factor	0.0004	3.77e-05	103.473	0.000	0.007	0.007
lagging_current_power_Factor	0.0009	3.93e-05	125.783	0.000	0.005	0.005
Omnibus:	4692.218	Durbin-Watson:	0.913			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	0.0005, 9.117			
Skew:	0.104	Prob(JB):	0.98			
Kurtosis:	0.452	Cond. No.:	2.54			

Log Transform OLS Regression Model - CO2

Result and Interpretation of Model

Metric	Original CO ₂	Log CO ₂	Δ (Improvement)
Adjusted R ²	0.8966	0.8983	+0.0017 (negligible)
BIC	-368,582.64	-370,589.95	↓ 2,007 points (✓)
Skewness	~+1.15	0.1	(✓) 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 R², 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("Original 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"Original kWh Model BIC: {bic_kwh_log:.2f}, Adjusted R-squared: {model_kwh_log.rsquared_adj:.4f}")
print(f"Original kWh BIC: 161259.72, Adjusted R-squared: 0.9110")
```

Result and Interpretation of Model

Model	Adjusted R ²	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

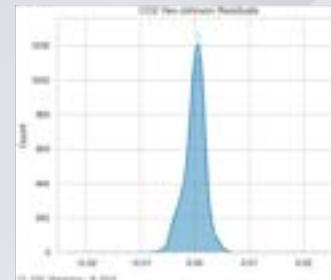
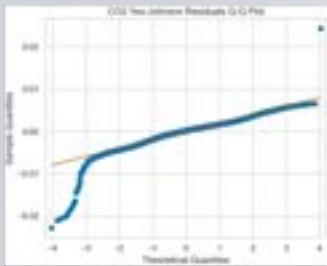
Since BIC penalizes for number of predictors; I assume that this drop isn't just due to added complexity. I believe it's a clear signal that the log transformation corrected skew and stabilized variance, aligning better with linear model assumptions.

```
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: 0.9110
```

Yeo-Johnson Model with CO₂

OLS Model (Yeo-Johnson)						
OLS Regression Results						
Dep. Variable:	y_1_001	R-squared:	0.954			
Model:	OLS	Adj. R-squared:	0.954			
Method:	Least Squares	F-statistic:	1,407.9e+06			
Date:	Mon, 10 Jun 2023	Prob (F-statistic):	0.00			
Time:	13:14:13	Log-Likelihood:	-1,487.9e+05			
N Observations:	15688	AIC:	-1,364.8e+05			
Df Residuals:	15684	BIC:	-1,367.4e+05			
Df Model:	4					
Covariance Type:	opg					
	(std. Err.)	t	P> t	[N. Obs]	R-sqrd	Adj. R-sqrd
const	8.0003	1.88e-05	584.768	0.000	0.955	0.954
temperature_2m (°C)	0.0007	1.23e-05	56.389	0.000	0.001	0.001
relative_humidity_2m (%)	-0.0003	1.24e-05	-17.966	0.000	-0.000	-0.000
Logging_Current_Basetime_Power_Norm	0.0036	1.31e-05	218.875	0.000	0.000	0.000
Logging_Current_Power_Factor	0.0048	1.48e-05	272.082	0.000	0.000	0.000
Logging_Current_Power_Factor	0.0054	1.54e-05	219.549	0.000	0.000	0.000
Deviance:	7285.428	Durbin-Watson:	0.458			
Prob>(Model):	0.000	Jarque-Bera (JB):	4.085e-06			
Score:	-0.028	Fprob (W):	0.00			
Kurtosis:	8.347	Gard. No.:	1.64			

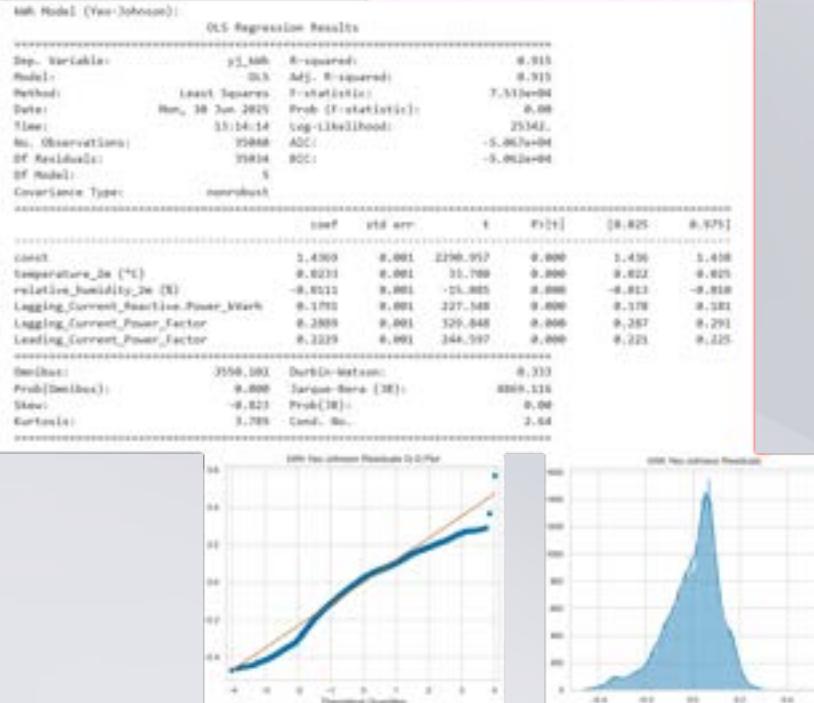


Result and Interpretation of Model

Metric	Value	Interpretation
Adjusted R ²	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 R²**
 - 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



Result and Interpretation of Model

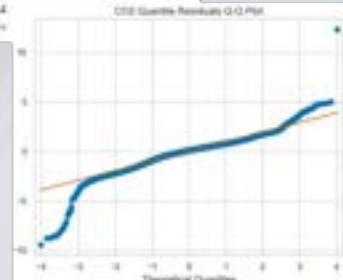
Metric	Value	Interpretation
Adjusted R ²	0.915	High explanatory power
BIC	-50,620	Substantial improvement from untransformed model
Skewness	-0.82	Left-skewed residuals with likely over suppression in low-usage cases
Kurtosis	3.79	Closer to normal with improved tail
Jarque-Bera	4,869.12	Residuals still 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 CO₂

OLS Regression Results											
Dep. Variable:	q1_CO2	R-squared:	0.897								
Model:	OLS	Adj. R-squared:	0.897 <th data-cs="4" data-kind="parent"></th> <th data-kind="ghost"></th> <th data-kind="ghost"></th> <th data-kind="ghost"></th>								
Method:	Least Squares	F-statistic:	5.13e+04								
Date:	Mon, 26 Jun 2023	Prob (F-statistic):	0.00								
Time:	13:14:15	Log-Likelihood:	-48759.								
N Observations:	35848	AIC:	9.745e+04								
Df Residuals:	35834	BIC:	9.798e+04								
Df Model:	5										
Covariance Type:	opg										
	tstat	std err	+/-	P> t	[9.0%	9.975]					
const	-2.7842	0.005	<2e-16	0.000	-2.795	-2.725					
Temperature_2m (°C)	0.3125	0.006	61.738	0.000	0.301	0.304					
relative_humidity_2m (%)	-0.0022	0.006	-25.117	0.000	-0.018	-0.000					
Lagging_Current_Active_Power_Watt	1.7517	0.007	287.582	0.000	1.741	1.766					
Lagging_Current_Power_Factor	0.8292	0.007	279.878	0.000	0.803	0.843					
Leading_Current_Power_Factor	1.7091	0.006	233.148	0.000	1.746	1.774					
Deviance:	5419.817	Durbin-Watson:	0.477								
Prob(Deviance):	0.000	Jarque-Bera (JB):	29138.317								
Score:	-0.448	Prob(JB):	0.00								
Kurtosis:	7.188	Cond. No.:	1.44								

Quantile transformation maps CO₂ 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 ²	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.

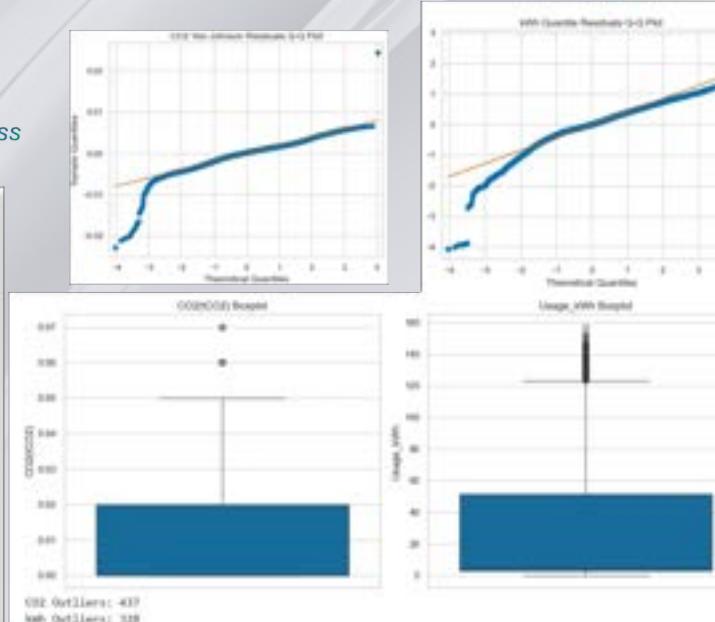
Checking for outliers for both CO₂ and Energy Usage kWh

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_kuh from previous
# analysis 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_kuh = remove_outliers(df['Usage_kWh'])

df_no_outliers = df[mask_co2 & mask_kuh].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_kuh'] = pt.fit_transform(df_no_outliers[['Usage_kWh']])
```

Metric	CO ₂ No Outliers	CO ₂ All Data	KWh No Outliers	KWh All Data
Adjusted R ²	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

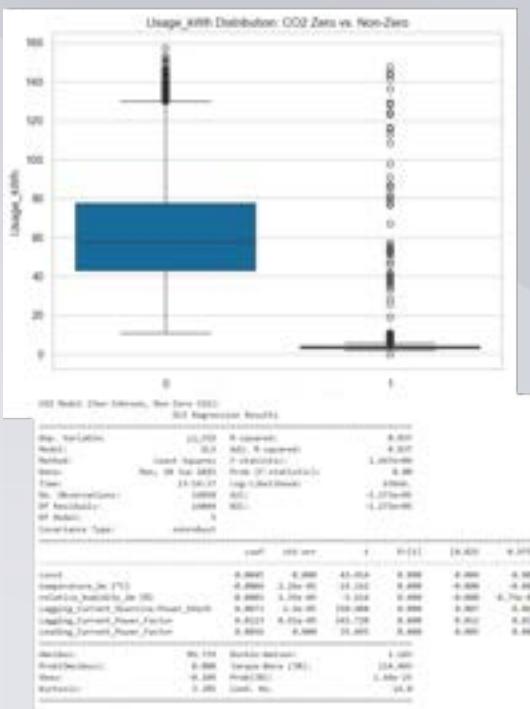
- BIC favors the full data, nudging toward keeping all points (even the quirky ones).
- Adjusted R-squared barely shifts, suggesting the models are resilient either way.
- Skewness and kurtosis steadily improve after trimming , especially in CO₂, which went from moderately skewed to more symmetric.

TAKEAWAY:

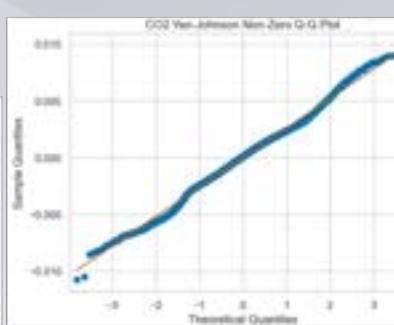
Outlier removal reduced tail risk (skewness, kurtosis) with minimal R² loss, but higher BIC suggests a trade-off in model parsimony.

CO₂ Non-Zero Model Analysis

BOXPLOT: CO₂ Non-Zero and Zero Comparison



 To better understand the structural behavior of CO₂ 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 R²) but hoped to improve distributional traits like skewness, kurtosis, and residual normality.



- **Adjusted R²** 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 CO₂ 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 sensure 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.

```

# Loading drivers for 100 countries or 500 data for multi-country country (2009-11-30). Download
# from: www.csi.org/csi.htm (or www.csi.org) provides the complete information (approximately 4,000
# data) downloaded approximately 200,000 in total bytes.

# Using 500 data. Model: 2002_plus_NBP <= y <= 2007 as new name CSD_NBP
# Define efficiency, using GMM production

# New statements, has vector or var model, import .mif

# Prepare data from csi.mif
dt_csi = dt_reduced[45,000 : 51,000][, -ncol(dt)]
csm = dt_csi[, names(dt_csi)[selected_hdp_columns(1)]]

# CSD_NBP
var_model = MIF2DF(dt_csi, wagegroup, csi_wagegroup, Amt = "C"
model = "CSD_NBP", Summary = "T")
print(var_model, summary = T)

# Estimated [summary]
# Estimated var_model, based on selected_hdp_columns(1) <- csm[, -csm$y], wage, polychorng(20, 2)
# print("CSD_NBP Forecast (20 years)") # print("CSD_NBP Forecast (20 years)", csm[, "y"], "CSD_NBP")
```

VAR Forecast Overview Next Slide



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.
 - $y_j_{\text{CO}_2}$, key predictors include Lagging_Current_Reactive.Power_kVarh (t-stat: 165.091) and Lagging_Current_Power_Factor (t-stat: 100.865).
 - For y_j_{kWh} , similar predictors are significant, with Leading_Current_Power_Factor showing a high t-stat (124.840).
 - **y_j_{kWh} equation working with $y_j_{\text{CO}_2}$'s history:** That L1 coefficient (0.5012, < 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 $y_j_{\text{CO}_2}$ suggest some oscillation.
 - Maybe seasonal behavior?

VAR Forecast Summary (24 Steps)

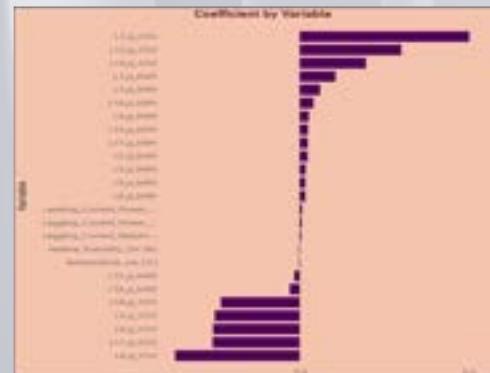
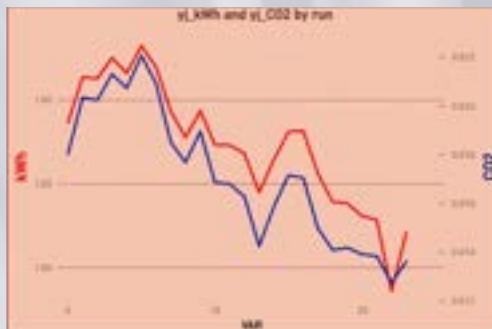
A 24-step-ahead forecast was generated using the last 24 observations of the endogenous variables and corresponding exogenous values. The forecast predicts y_{j_CO2} and y_{j_kWh} over a 24-hour period, reflecting expected patterns in emissions and electricity use.

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 for CO2 emissions and energy usage high level code used for capturing the historical effects of energy usage and consumption. The VAR model uses and records the historical regression history to capture lagged causalities and shifts, prior to fitting coefficients to capture causality (e.g., high load will increase CO2 to reflect user activity and portfolio demand).  
from statsmodels.tsa.api import VAR  
from statsmodels.tsa.stattools import acf  
from statsmodels.tsa.stattools import grangercausalitytests  
  
# Process data  
df_var = df[['consumed_kWh', 'yj_CO2']].head(1000)  
exog = df[['consumed_kWh', 'yj_CO2']].tail(1000)  
  
# fit VAR  
var_model = VAR(df_var, endog=df['yj_CO2'])  
print("Fitted Model Summary")  
print(var_model.summary())  
  
# fit exogenous  
print("Fitted lag coefficients yj_CO2 vs consumed_kWh")  
print(var_model.coefs[0])  
print("Fitted lag coefficients yj_CO2 vs consumed_kWh")  
print(var_model.coefs[1])  
  
# Diagnostic statistics  
print("Fitted residuals P-value: ", var_model.pvalues)  
print("Fitted residuals Q-value: ", var_model.qvalues)  
print("Fitted residuals Durbin-Watson: ", var_model.durbin_watson)  
  
# JLR plot  
var_model.irf(24).plot(orth=True)  
plt.show()  
  
# Test MSE (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=100, exog=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 CO₂ 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.
- CO₂ 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 CO₂ 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 CO₂ output. The demand threshold (kWh_high) amplifies this response, making it an operational pivot point.

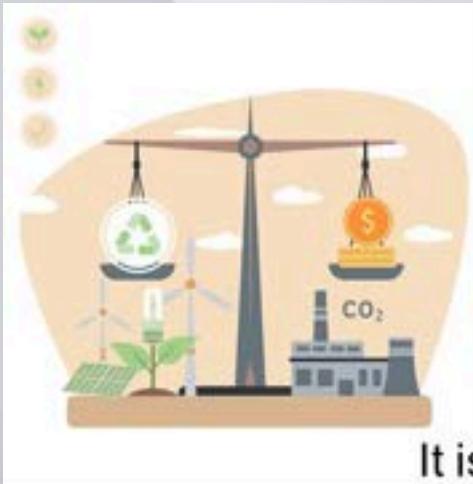
CO₂ Emission Highlights

- kWh_high: Negative impact (-0.00173, p < 0.001) with higher usage links to lower marginal emissions.
- yj_CO₂ 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_CO₂: 0.000007, really looking good!



Interpretation

CO₂ 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?"

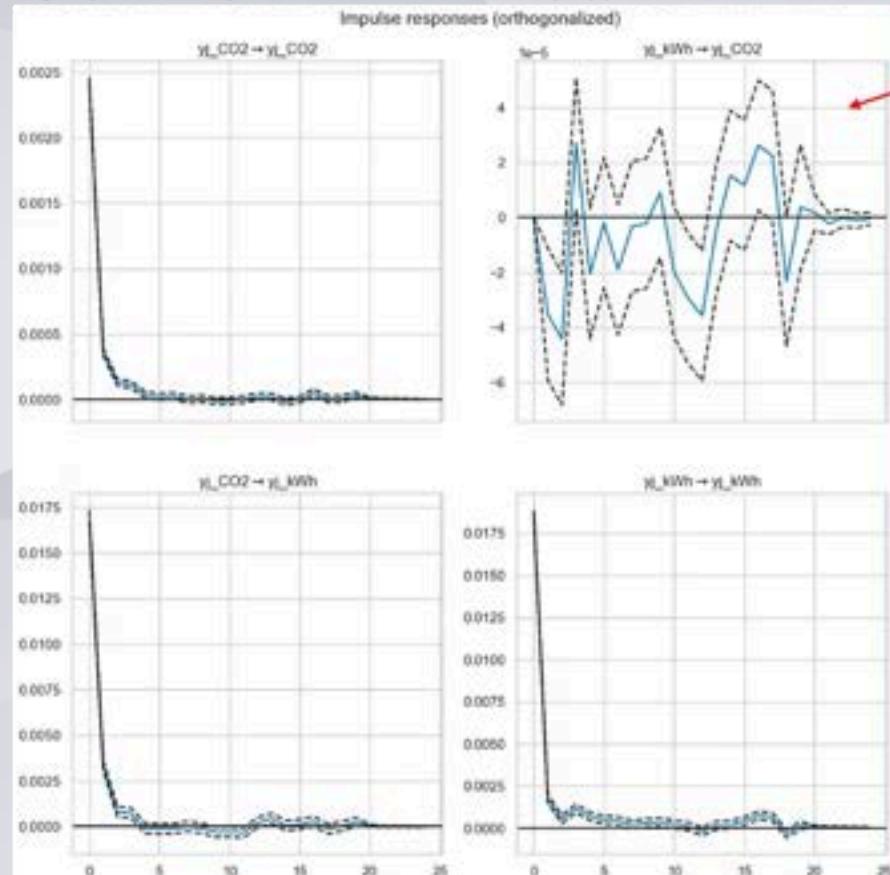
Technical Summary Table

Test	Null Hypothesis	Result	p-value	Interpretation
$kWh \rightarrow CO_2$	kWh does not Granger-cause CO_2	✗ Rejected	0	Energy usage helps forecast emissions.
$CO_2 \rightarrow kWh$	CO_2 does not 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 \rightarrow Response	Pattern	Interpretation
$y_{j,CO_2} \rightarrow y_{j,CO_2}$	Small initial bump (~0.0025), fades fast	Emissions have short memory; autoregressive persistence is limited.
$y_{j,CO_2} \rightarrow y_{j,kWh}$	Quick drop (~0.0175), stabilizes	Energy reacts swiftly to CO ₂ spikes — suggests fast operational feedback.
$y_{j,kWh} \rightarrow y_{j,kWh}$	Immediate decline, small oscillations	Self-stabilizing; supports high-frequency autoregression in usage.
$y_{j,kWh} \rightarrow y_{j,CO_2}$	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 y_{j,CO_2} and $y_{j,kWh}$

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

```
Johansen's Test
We compute Johansen (1988) for cointegration.
H0: H0: no cointegration (null)
H1: H1: there is at least one cointegrating relationship
Note: when you have two or more time series that's available in the long run but might contain some
short-term error, there's a cointegrating relationship.
P-value: Q: Above contains a critical value of rank 1, see H0P.
H0 null, H1: H0H1 is sufficient.

Matrix for co-integrating long-run (1988 and 2001) and of cointegration
which might go up/down daily (co-integration)
a highly correlated perfectly linear
that long-term they move together
that T-invariant differences due to seasonal relationships

from statsmodels.tsa.cointintegration import coint_johansen
import pandas as pd

# download df_johansen
df_johansen = pd.read_csv('yj_CO2_yj_kWh.csv')

# Johansen test
johansen = coint_johansen(df_johansen['yj_CO2'], df_johansen['yj_kWh'], 1)
print(johansen.concise_pvalues[1])
print(johansen.concise_pvalues[2])
print(johansen.concise_pvalues[3])
print(johansen.concise_pvalues[4])
```

Rank	Trace Statistic	Critical Value (5%)	Eigen Statistic	Critical Value (5%)	Conclusion
0	5476.76	15.49	3837.86	14.26	Reject H_0
1	1638.9	3.84	1638.9	3.84	Reject H_0

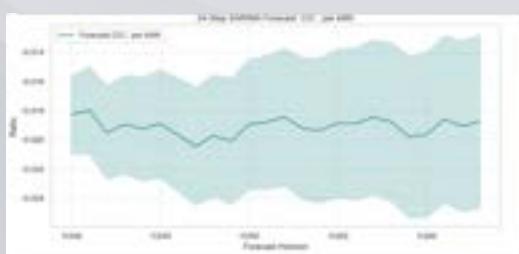
Interpretation

- 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^2 for CO_2 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 – CO₂ per kWh

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



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 ~-0.24 — striking signal regularity.

- Even though CO₂ per kWh had low R² 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 option for kWh to capture non-linearity)
df_nonzero['kWh_high'] = df_nonzero['Usage_kWh'] > 20).astype(int)

X = x.add_constant(df_nonzero[['Usage_kWh']] > 20).as_matrix()
y = df_nonzero['CO2_per_kWh']

model_CO2_per_kWh = sm.OLS(y, X).fit()
print("model_CO2_per_kWh with kWh_high:")
print(model_CO2_per_kWh.summary())
```

CO2_per_kWh with kWh_high						
OLS Regression Results						
Model:	CO2_per_kWh	OLS	R-squared:	0.814		
Method:	Least Squares		F-statistic:	0.020000		
Date:	Mon, 07 Nov 2011		Prob (F-statistic):	0.40		
Time:	19:02:19		R-squared:	0.014		
Df Residuals:	1484		Adj R-squared:	0.014		
Df Model:	4		F-statistic:	0.475e+00		
Convergence Type:	converged					
	const	std err	t	AIC	BIC	LL
	-0.0012	0.000	-0.001	0.000	0.000	0.000
Intercept:	0.0000	0.000	0.000	0.000	0.000	0.000
relatives_homedist_km (O)	-0.0000	0.000	-0.000	0.000	0.000	0.000
logging_currents_homedist_km	0.0000	0.000	0.000	0.000	0.000	0.000
relative_homedist_km_sq	0.0000	0.000	0.000	0.000	0.000	0.000
logging_currents_homedist_sq	0.0000	0.000	0.000	0.000	0.000	0.000
relative_homedist_sq	0.0000	0.000	0.000	0.000	0.000	0.000
logging_currents_homedist_sq_sq	0.0000	0.000	0.000	0.000	0.000	0.000
relatives_homedist_sq_sq	0.0000	0.000	0.000	0.000	0.000	0.000
Std. Err:	0.000	0.000	0.000	0.000	0.000	0.000
Residual Std. Dev:	0.000	0.000	0.000	0.000	0.000	0.000
Residuals:	0.000	0.000	0.000	0.000	0.000	0.000
Warning: Standard Errors assume that the covariance matrix of the errors is correctly specified.						
(2) The condition number is large, 3.3e+07. This might indicate that there are strong multicollinearity or other numerical problems.						

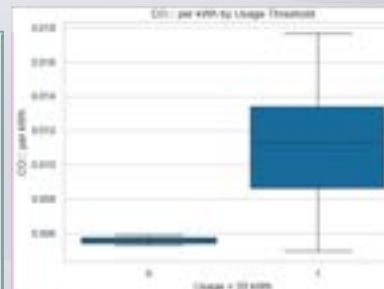
What This Nonlinear Model Reveals

- **R² = 0.814** is a dramatic leap from initial CO₂_per_kWh regression (R² ~ 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. CO₂ 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

CO₂ efficiency doesn't follow a straight line

Get policy digest manifest
-Writing a cache option with pre-verified capture items to the same-saving path for utilization
-Utilizing a file for cache pre-charge, file
From policy digest manifest

```

# Create a matrix for these NPs
n_lines = matrix(NA, nrow=NP, ncol=NP)
for (i in 1:NP) {
  for (j in 1:NP) {
    if (i == j) {
      n_lines[i,j] = 1
    } else {
      n_lines[i,j] = sum((NP$NP[i] %in% NP$NP[j]) | (NP$NP[j] %in% NP$NP[i]))
    }
  }
}
# Create a matrix for the number of tokens selected but not used
t_lines = matrix(NA, nrow=NP, ncol=NP)
for (i in 1:NP) {
  for (j in 1:NP) {
    if (i == j) {
      t_lines[i,j] = 1
    } else {
      t_lines[i,j] = sum((NP$tokens[i] %in% NP$tokens[j]) | (NP$tokens[j] %in% NP$tokens[i]))
    }
  }
}
# Create a matrix for the number of tokens used
u_lines = matrix(NA, nrow=NP, ncol=NP)
for (i in 1:NP) {
  for (j in 1:NP) {
    if (i == j) {
      u_lines[i,j] = 1
    } else {
      u_lines[i,j] = sum((NP$tokens[i] %in% NP$tokens[j]) & (NP$NP[i] %in% NP$NP[j]))
    }
  }
}

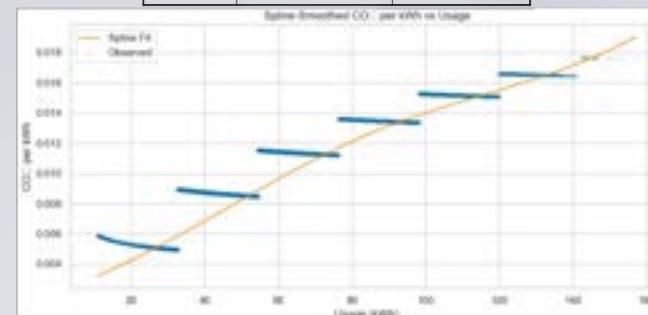
```

```
print(saville_national_summary())
```

Notes:
1) Standard Survey assumes that the conversion matrix of the account is correctly specified.
2) The coefficient number is 2.00048. This might indicate that there are

Spline regression revealed soft structural curvature in CO₂-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
$R^2 = 0.908$	Excellent fit	Found structure?
Durbin-Watson ~1.676	Mild residual autocorrelation	Time-based dependencies matter
Curve shape	Gentle rise, mild bends	Suggests controlled efficiency scaling
Scatter points	Slight curvature in actuals	Affirms nonlinearity in observed data

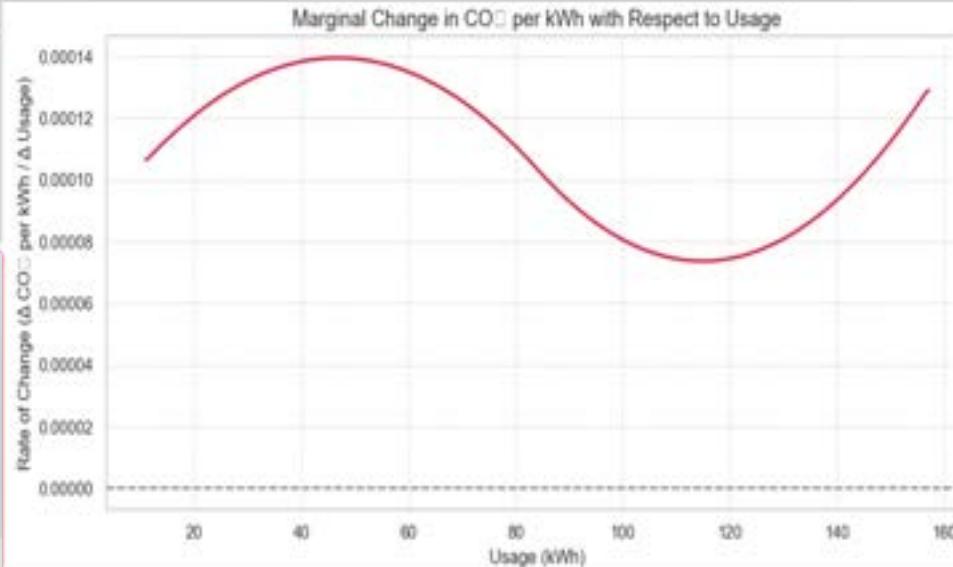


Spline-Augmented Model on Usage_kWh

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 ~ 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 (≈ 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_coor = df_johansen[['L1_CO2', 'L1_kWh']].dropna()
df_kwh = df_johansen.dropna().reset_index().columns[11:16].dropna()
df_coor['Date'] = df_coor.index.to_datetime()
df_kwh['Date'] = df_kwh.index.to_datetime()

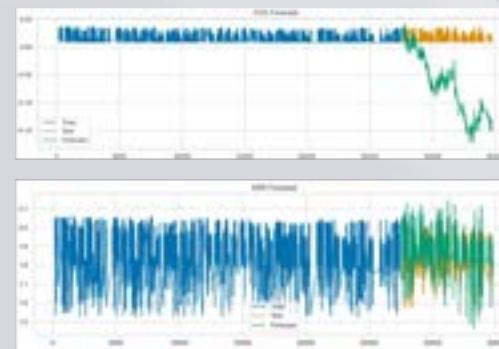
# Train-test split
train_size = int(len(df_coor) * 0.8)
train_coor, test_coor = df_coor[:train_size], df_coor[train_size:]
train_kwh, test_kwh = df_kwh[:train_size], df_kwh[train_size:]

# Fit VECM
coint_mdl = VECM(coor, kwh, exog_lags=1, exog_kwh_lags=1, exog_coor_lags=1, exog_coor_kwh_lags=1)
print("Output Model Summary:")
print(coint_mdl)

# Forecast
yhat_coor = coint_mdl.predict(1000, exog_coor_lags=1, exog_kwh_lags=1)
yhat_kwh = coint_mdl.predict(1000, exog_coor_lags=1, exog_kwh_lags=1, exog_coor_kwh_lags=1)

# Test MSE
test_coor['mean_squared_error'] = np.mean((test_coor['CO2'].values - yhat_coor)**2)
test_kwh['mean_squared_error'] = np.mean((test_kwh['kWh'].values - yhat_kwh)**2)
print("Test MSE: CO2 = {} | kWh = {}".format(test_coor['mean_squared_error'], test_kwh['mean_squared_error']))
print("Forecast MSE: CO2 = {} | kWh = {}".format(yhat_coor[-1], yhat_kwh[-1]))

```



VECM MODEL



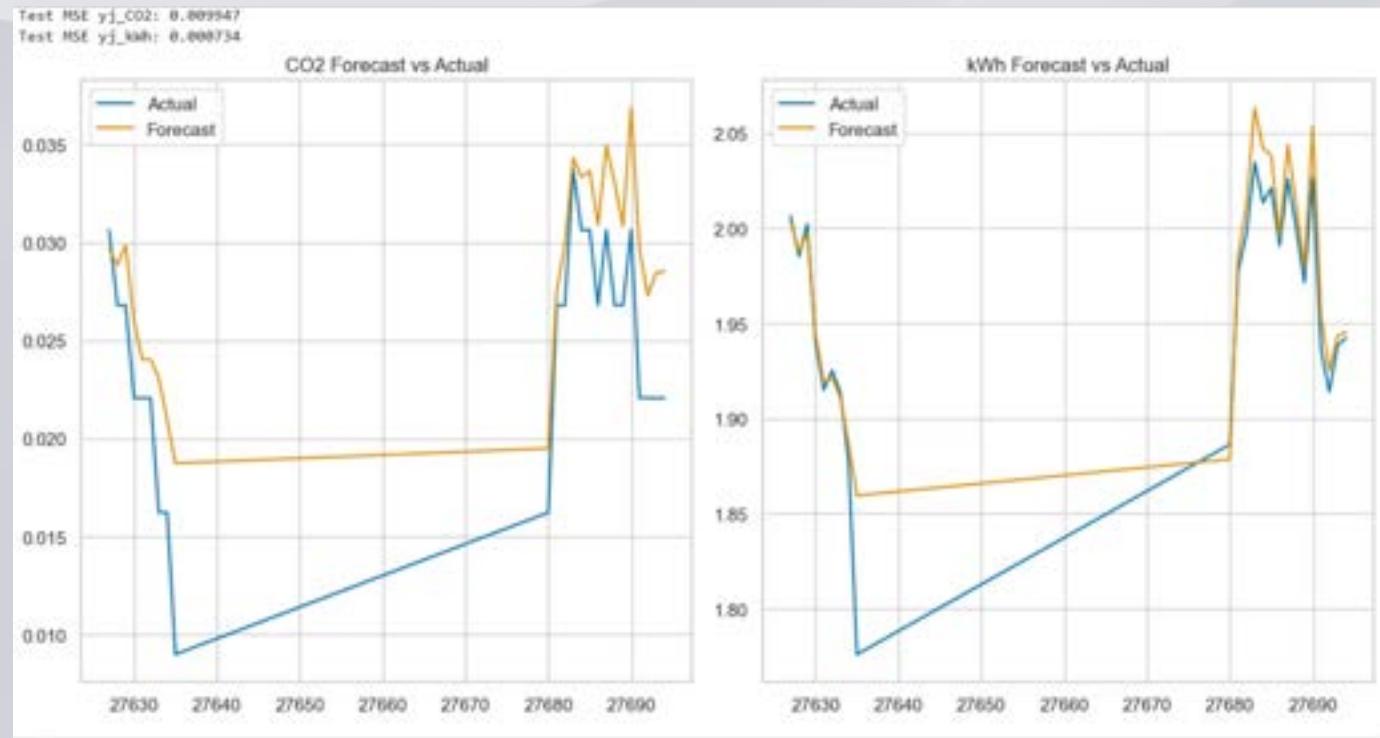
Returning to the vector error correction model (VECM) which is an application of vector autoregression (VAR) that incorporates error correction terms to capture long-run relationships with **cointegrated** variables. The VECM can also estimate both short-run and long-run coefficients. Using my past VAR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

	coef	std. err.	t	p-value	(0.05)	0.5%
Coef. table (Model 1)						
coor[0]	-0.0001	0.0001	-0.704	0.480	-0.101	
coor[1]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[2]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[3]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[4]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[5]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[6]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[7]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[8]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[9]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[10]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[11]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[12]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[13]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[14]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[15]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[16]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[17]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[18]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[19]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[20]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[21]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[22]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[23]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[24]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[25]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[26]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[27]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[28]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[29]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[30]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[31]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[32]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[33]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[34]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[35]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[36]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[37]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[38]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[39]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[40]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[41]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[42]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[43]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[44]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[45]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[46]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[47]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[48]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[49]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[50]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[51]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[52]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[53]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[54]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[55]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[56]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[57]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[58]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[59]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[60]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[61]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[62]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[63]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[64]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[65]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[66]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[67]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[68]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[69]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[70]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[71]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[72]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[73]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[74]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[75]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[76]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[77]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[78]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[79]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[80]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[81]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[82]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[83]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[84]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[85]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[86]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[87]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[88]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[89]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[90]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[91]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[92]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[93]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[94]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[95]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[96]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[97]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[98]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[99]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[100]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[101]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[102]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[103]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[104]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[105]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[106]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[107]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[108]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[109]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[110]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[111]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[112]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[113]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[114]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[115]	0.452e-05	0.154e-05	2.945	0.003	0.050e-05	-0.200e-05
coor[116]	0.452e-05	0.154e-05	2.945	0.003	0.0	

Better Look of Forecast vs. Actual for VCEM Model

Forecast vs. Actual: Detecting Equilibrium Drift and Recovery

Test MSE y_{j_CO2} : 0.009947
Test MSE y_{j_kWh} : 0.000734



Re-Running ARIMA

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



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 ($D=1$, $s=\text{seasonal period}:124$ 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.

Looking at the df_nonzero['yj_CO2'] Distribution



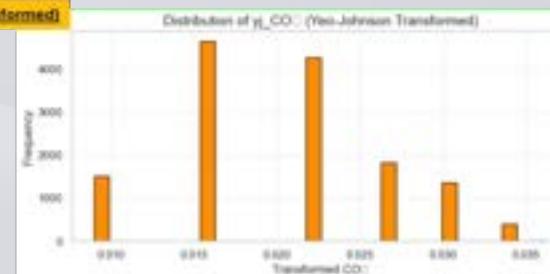
Before diving into more complex models, it's important to understand the shape and behavior of my transformed variable. Exploring the distribution of yj_CO₂ ensures that the assumptions of normality, variance stability, and scale suitability are met. This step acts as a diagnostic check and clearing 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=10, color='darkorange', edgecolor='black')  
plt.title("Distribution of yj_CO2 (Yeo-Johnson Transformed)", fontweight='bold')  
plt.xlabel("Transformed CO2")  
plt.ylabel("Frequency")  
plt.grid(True, linestyle="--", alpha=0.1)  
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_CO₂ (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.)

```
RandomForestRegressor(n_estimators=100, max_depth=5, min_samples_leaf=1, min_weight_fraction_leaf=0.001, max_features='auto', max_leaf_nodes=None, bootstrap=True, oob_score=False, n_jobs=-1, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)

RandomForestRegressor(max_depth=5, max_features='auto', max_leaf_nodes=None, min_samples_leaf=1, min_weight_fraction_leaf=0.001, n_estimators=100, n_jobs=-1, oob_score=False, random_state=None, verbose=0, warm_start=False)
```

RANDOM FOREST MODEL

With main factors from previous models

Okay, I am done! How could I improve 😊

```
fit Score: 0.9999999991728026  
RMSE: 5.053941883082775e-06
```

Feature Importance:

	feature	importance
5	Usage_kWh	1.0
8	temperature_2m ("C")	0.0
1	relative_humidity_2m (%)	0.0
2	Logging_Current_Reactive_Power_kVarh	0.0
3	Logging_Current_Power_Factor	0.0
4	Leading_Current_Power_Factor	0.0

```
Cross-validation scores: [0.99972452 0.99990276 0.99999109 0.99999965 0.99999457]  
Mean CV score: 0.999922519832897
```

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 (CO₂ without Usage_kWh)

RF Model 2

```
#Not surprised that kWh was the dominating factor,
#Assessing two models - one model for CO2 without kWh and CO2 without not
#RandomForest Model CO2 without kWh

#Correlated variables - remove usage_kWh - Only note to check later()
X_cold = df['measured']
    'temperature_2m (°C)', 
    'relative_humidity_2m (%)',
    'Lagging_Current_Reactive_Power_kVarh',
    'Lagging_Current_Power_Factor',
    'Leading_Current_Power_Factor'
    ]) # remove: no usage_kWh here
y_cold = df['measured'].values[1:1000]

#Split data
X_train, X_test, y_train, y_test = train_test_split(X_cold, y_cold, 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 assess
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_cold.columns, # Using X_cold 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:

- | Rank | 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

CO₂ without Usage_kWh

- R² = 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.

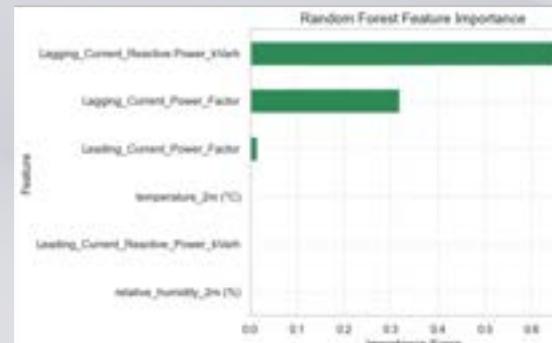
RANDOM FOREST (Usage_kWh without co₂)

RF Model 3

R2 Score: 0.997504913036296
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
```



Overview of Model Random Forest Model 3

Usage_kWh without CO₂

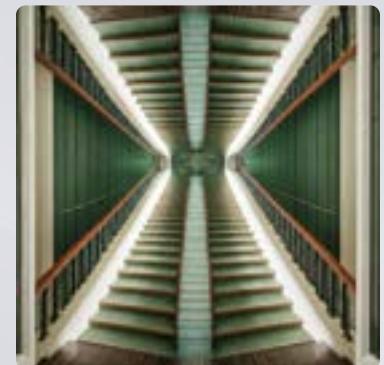
- Similar pattern to the CO₂-without-kWh model.
 - Again, **Reactive Power and Power Factors dominate**, suggesting strong internal correlations between electrical dynamics and usage.
 - Temperature and humidity again make only a slight appearance.

Dual Models: CO₂ Emissions vs. Energy Usage

	CO ₂ Model (without kWh)	kWh Model (without CO ₂)
R ² 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 CO₂ 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?

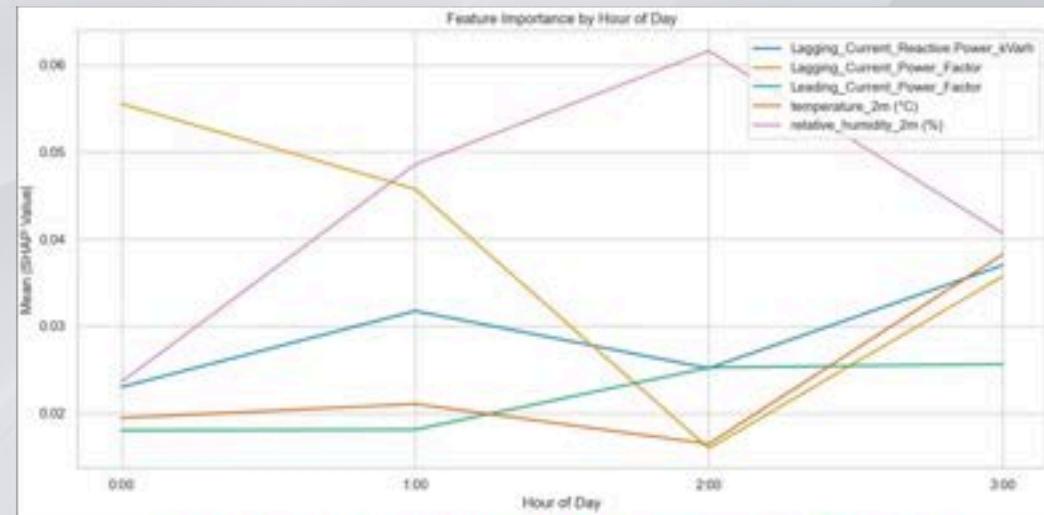


SHAP (SHapley Additive exPlanations)

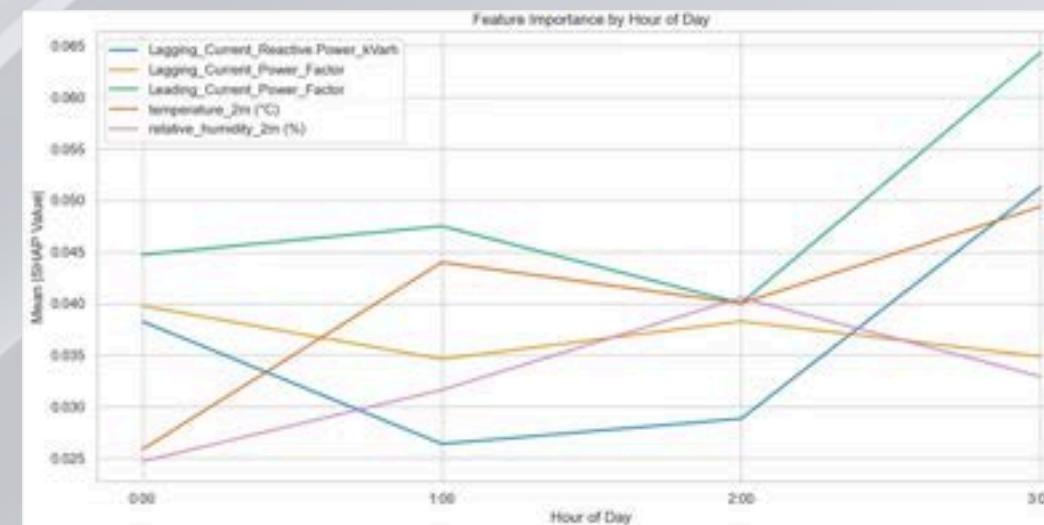
“SHAP explains how each feature contributes to a model’s prediction by assigning an impact score to every variable for every individual prediction. It’s like giving your data a voice in how the model made its decision.”

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.

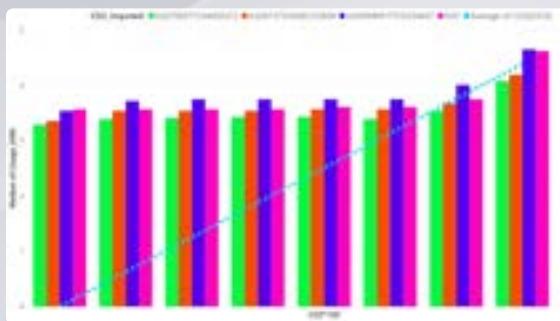
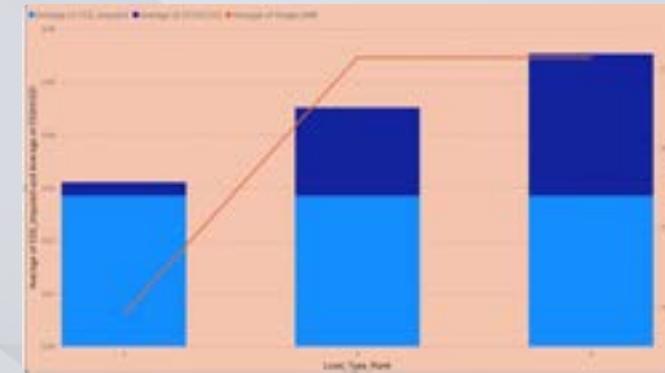
Reconsidering Imputed CO₂ in Predictive Modeling

While early modeling benefited from imputed CO₂ 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 CO₂.

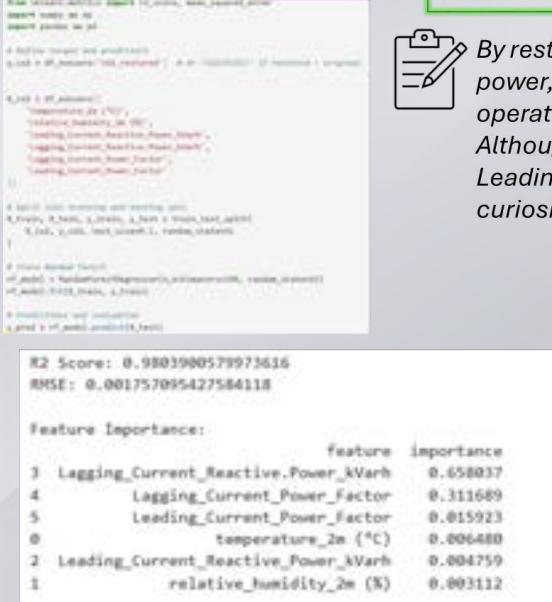
Rank 2 approached a 50/50 split, still skewed toward imputation.

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



A complementary bar chart of imputed-only CO₂ demonstrated another red flag: values didn't increase in a linear fashion. Unlike observed CO₂, imputed levels didn't correspond naturally with operational intensity, creating unexpected elevation patterns where zeros weren't the lowest in the group.

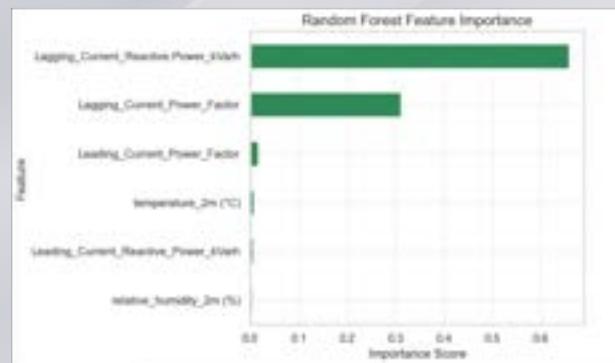
Random Forest Model With Restored CO₂



A small icon of a clipboard with a checklist.

By restoring the original CO₂ 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.



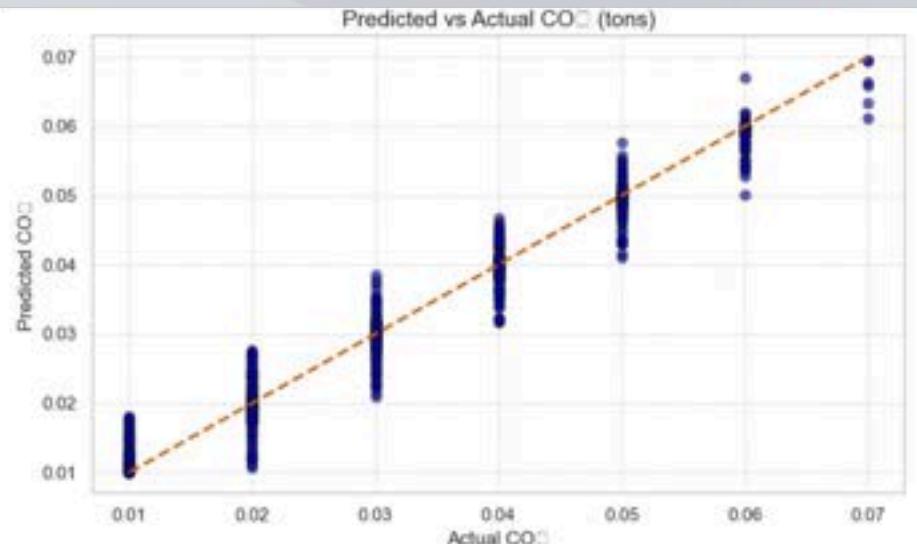
OVERVIEW ON NEXT SLIDE



Random Forest Model With Restored CO₂

```
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:

- **R² 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 CO₂.**

Takeaways:

- Reactive energy and power factor still dominate CO₂ 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 CO₂ 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 CO₂ groupings which is another signal that the model is capturing structured behavior in emissions.

Partial Dependence Insights of CO₂ Random Forest

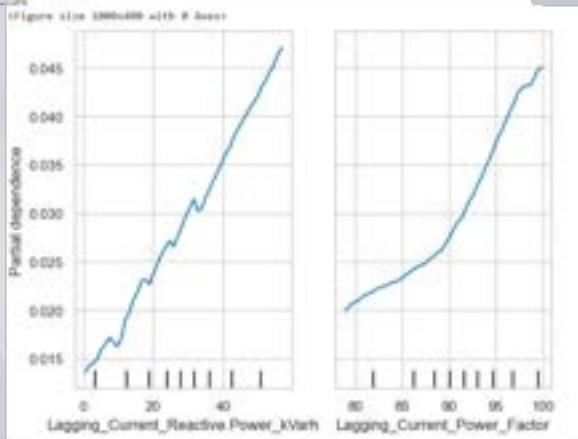
```
from sklearn.inspection import PartialDependenceDisplay
import matplotlib.pyplot as plt

Features = ['Lagging_Current_Reactive_Power_kVarh', 'Lagging_Current_Power_Factor']

# Set the figure size before importing the module
plt.figure(figsize=(10, 4))

# Create and plot the partial dependence
display = PartialDependenceDisplay.from_estimator(
    rf_model, X_test, Features, kind='average', grid_resolution=100,
    feature_namesR_idx.columns
)

display.plot()
plt.tight_layout()
```



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 CO₂ emissions escalate.
- May reflect systems shifting into less efficient states or triggering compensatory responses after certain load levels.

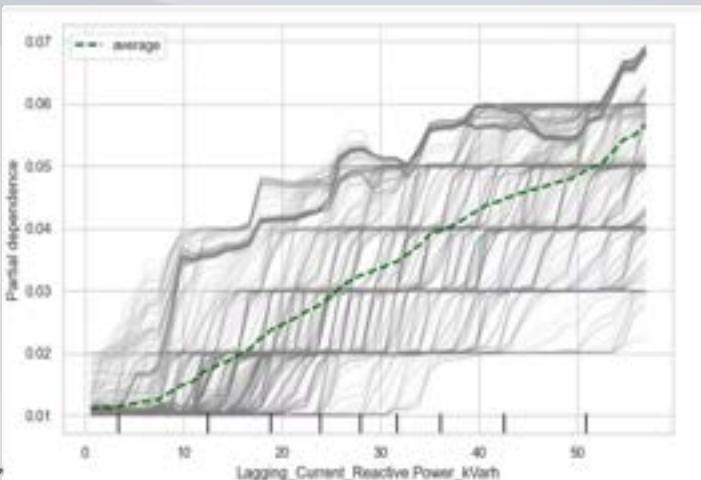
Lagging_Current_Power_Factor

- 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 CO₂.
 - This is potentially due to load balancing or control system effects.

ICE PLOT OF LAGGING Current Reactive Power KVarh

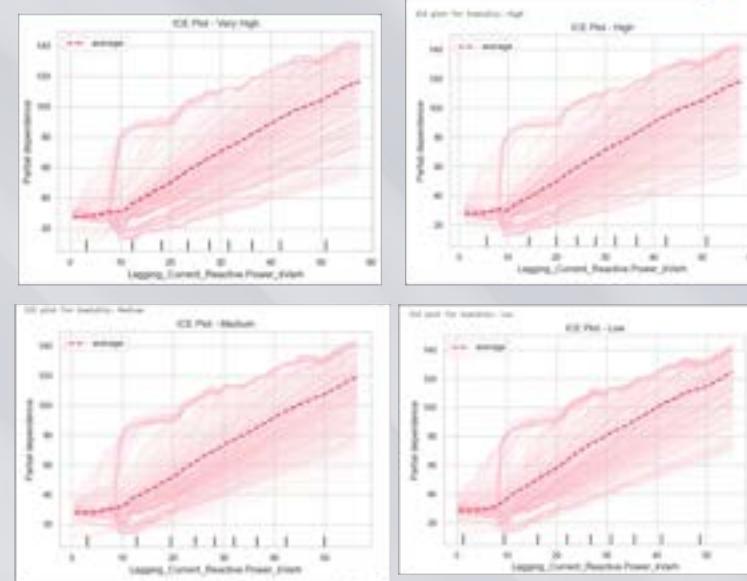
```
What ICE plots offer - While the MPD line shows the average effect of a feature (like a weighted mean),  
ICE plots reveal how different samples react to changing just that one feature.  
It's like multiple person's stories mixed until the summary:  
Printing IP for me has features in the model: Lagging_Current_Reactive_Power_KVarh  
From vision, inspection import PartialDependencePloting  
import matplotlib.pyplot as plt  
  
# Create figure with desired size FigSize  
fig, ax = plt.subplots(figsize=(8, 12))  
  
# Individual curves for each sample = MPD curves  
PartialDependencePloting.from_estimator(  
    pf_model, X_test,  
    feature='Lagging_Current_Reactive_Power_KVarh',  
    n_jobs=-1, n_reps=100, n_folds=5,  
    int_size=(1000, 600), color='gray'), # A matrix 100 rows  
    pd_line_low_value=1, 'background': 'white', 'grid_color': 'black',  
    feature_name='Lagging_Current_Reactive_Power_KVarh', # Make sure this variable is defined  
    grid_resolution=10,  
    axes = fig) # Pass the axes object instead of Figure  
)  
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 CO₂ 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.

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



- **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**.

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 CO₂) COMPARISON: Adding and Reducing Variables

```

# Model for full-cost additional methods
class CPMMethods:
    "Initialization for CPM"
    "Initialization, based on BOM"
    "Logging Current, Specific Power Factor"
    "Logging Current, Power Factor"
    "Logging Current, Power Factor"
    "Logging Current, Specific Power Factor"
    "Logging Current, Specific Power Factor"
    "Logging Current, Specific Power Factor"

# Starting user model, setup as factory, via factory of additional layer of nesting
# everything
# A specific model
CPMMethods(1,2000, 1,2000 + 10000, 100, 10000, 1, 100, 10000, 1, 10000, 10000)

# Create and create model
# of model = RandomizedRegressionModel, set parameters, use Asymptotic, random_element()
# of model = PLSR, k=5000

# Run predictions and evaluate
# predict = fit_random_forest(X, y)
# predict = M.lasso(), fit_random_forest(X, y)
# predict = MRL, fit_random_forest(X, y)
# predict = MRL, fit_random_forest(X, y, n_estimators=10000)

# Feature importance
importance = fit_random_forest(X, y).feature_importances_
# Feature = X[fit_random_forest(X, y).feature_importances_.argsort()]
# importance = fit_random_forest(X, y).feature_importances_

```

R2 Score: 0.9983005599560189
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	0.814047
9	temperature_2m (°C)	0.002462
5	Leading_Current_Reactive_Power_kVarh	0.001795
1	relative humidity_2m (%)	0.000578

R2 Score: 0.96956578968678
RMSE: 0.17464776994887217

Feature Importance

	feature	importance
0	Lagging_Current_Reactive_Power_VArh	0.63964
1	Lagging_Current_Power_Factor	0.29115
2	Leading_Current_Power_Factor	0.06919

Cross-validation scores: [0.94688156 0.95683223 0.96954724 0.97113622 0.96314634]
Mean CV score: 0.961348718241472

RANDOM FOREST (CO₂ without Usage_kWh) Streamlining Variables

After seeing how well the streamlined model performed for Energy Usage (kWh), I couldn't resist revisiting the CO₂ 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_val = df['measure']
x_val = x_val[['Lagging_Current_Reactive_Power_Mvarh',
               'Lagging_Current_Power_Factor',
               'Leading_Current_Power_Factor']]
# Removed Temperature, humidity, and Leading_Current_Reactive_Power_Mvarh
x_val = x_val[['measure', 'p1_000']]

# Split data
x_train, x_test, y_train, y_test = train_test_split(x_val, y_val, test_size=0.1, 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_rf = pd.DataFrame()
importance_rf['X_val'] = x_val.columns
importance_rf['Importance'] = rf_model.feature_importances_
print('Feature importance')
print(importance_rf.sort_values('Importance', ascending=False))

# Cross-validation
scores = cross_val_score(rf_model, x_val, y_val, cv=5)
print("Cross-validation scores: ", scores)
print("Mean CV score: ", scores.mean())
```

```
R2 Score: 0.969565780686782
RMSE: 0.1746476994087217

Feature Importance:
          feature      importance
0  Lagging_Current_Reactive_Power_Mvarh  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 CO₂) 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 ² Score	RMSE	Top Features
Full Model (6 features)	0.9983	1.0966	Lagging Reactive Power, Lagging Power Factor
Streamlined (3 features)	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 (CO₂ without Usage_kWh) Streamlining Variables

Comparative Summary

	Full Model (6 features)	Streamlined Model (3 features)
R² Score	0.9804	0.9696
RMSE	0.00176	0.17465
Top Feature	Lagging Reactive Power (65.8%)	Lagging Reactive Power (63.96%)
Secondary Features	Lagging PF (31.2%), Leading PF (1.6%)	Lagging PF (29.12%), Leading PF (6.92%)
Environmental Variables	Present (humidity, temp, etc.) — minimal impact	Excluded — focus purely on core electrical behavior
Added Variable	Leading Reactive Power (small but above humidity)	Not included
Cross-Validation Mean	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



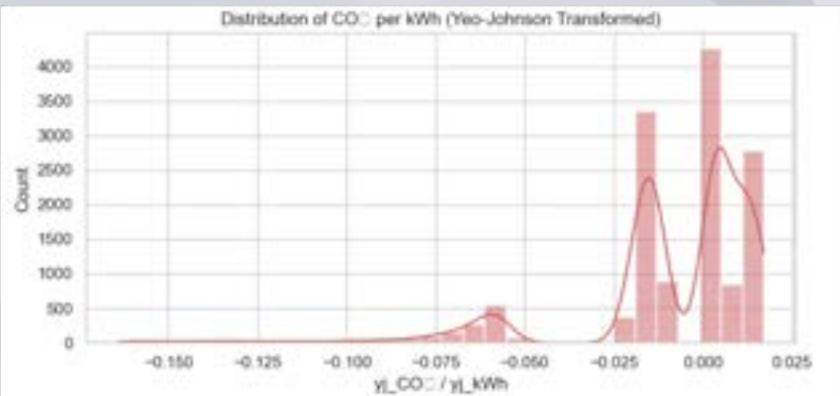
With time series and seasonal structure showing strong influence across models, I revisited the VECM forecast to observe how emissions and energy usage unfold over time. The following plots reflect not only directional alignment but also systemic behavior across hours and operational states, underscoring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

	coeff	std. err.	t	P> t	[<0.05]	>0.05]
co2	-0.000	0.100	-0.000	0.427	0.200	
co2p1	0.000	0.000	0.000	0.000	0.000	
co2p2	0.000	0.000	0.000	0.000	0.000	
co2p3	0.000	0.000	0.000	0.000	0.000	
co2p4	0.000	0.000	0.000	0.000	0.000	
co2p5	0.000	0.000	0.000	0.000	0.000	
co2p6	0.000	0.000	0.000	0.000	0.000	
co2p7	0.000	0.000	0.000	0.000	0.000	
co2p8	0.000	0.000	0.000	0.000	0.000	
co2p9	0.000	0.000	0.000	0.000	0.000	
co2p10	0.000	0.000	0.000	0.000	0.000	
co2p11	0.000	0.000	0.000	0.000	0.000	
co2p12	0.000	0.000	0.000	0.000	0.000	
co2p13	0.000	0.000	0.000	0.000	0.000	
co2p14	0.000	0.000	0.000	0.000	0.000	
co2p15	0.000	0.000	0.000	0.000	0.000	
co2p16	0.000	0.000	0.000	0.000	0.000	
co2p17	0.000	0.000	0.000	0.000	0.000	
co2p18	0.000	0.000	0.000	0.000	0.000	
co2p19	0.000	0.000	0.000	0.000	0.000	
co2p20	0.000	0.000	0.000	0.000	0.000	
co2p21	0.000	0.000	0.000	0.000	0.000	
co2p22	0.000	0.000	0.000	0.000	0.000	
co2p23	0.000	0.000	0.000	0.000	0.000	
co2p24	0.000	0.000	0.000	0.000	0.000	
co2p25	0.000	0.000	0.000	0.000	0.000	
co2p26	0.000	0.000	0.000	0.000	0.000	
co2p27	0.000	0.000	0.000	0.000	0.000	
co2p28	0.000	0.000	0.000	0.000	0.000	
co2p29	0.000	0.000	0.000	0.000	0.000	
co2p30	0.000	0.000	0.000	0.000	0.000	
co2p31	0.000	0.000	0.000	0.000	0.000	
co2p32	0.000	0.000	0.000	0.000	0.000	
co2p33	0.000	0.000	0.000	0.000	0.000	
co2p34	0.000	0.000	0.000	0.000	0.000	
co2p35	0.000	0.000	0.000	0.000	0.000	
co2p36	0.000	0.000	0.000	0.000	0.000	
co2p37	0.000	0.000	0.000	0.000	0.000	
co2p38	0.000	0.000	0.000	0.000	0.000	
co2p39	0.000	0.000	0.000	0.000	0.000	
co2p40	0.000	0.000	0.000	0.000	0.000	
co2p41	0.000	0.000	0.000	0.000	0.000	
co2p42	0.000	0.000	0.000	0.000	0.000	
co2p43	0.000	0.000	0.000	0.000	0.000	
co2p44	0.000	0.000	0.000	0.000	0.000	
co2p45	0.000	0.000	0.000	0.000	0.000	
co2p46	0.000	0.000	0.000	0.000	0.000	
co2p47	0.000	0.000	0.000	0.000	0.000	
co2p48	0.000	0.000	0.000	0.000	0.000	
co2p49	0.000	0.000	0.000	0.000	0.000	
co2p50	0.000	0.000	0.000	0.000	0.000	
co2p51	0.000	0.000	0.000	0.000	0.000	
co2p52	0.000	0.000	0.000	0.000	0.000	
co2p53	0.000	0.000	0.000	0.000	0.000	
co2p54	0.000	0.000	0.000	0.000	0.000	
co2p55	0.000	0.000	0.000	0.000	0.000	
co2p56	0.000	0.000	0.000	0.000	0.000	
co2p57	0.000	0.000	0.000	0.000	0.000	
co2p58	0.000	0.000	0.000	0.000	0.000	
co2p59	0.000	0.000	0.000	0.000	0.000	
co2p60	0.000	0.000	0.000	0.000	0.000	
co2p61	0.000	0.000	0.000	0.000	0.000	
co2p62	0.000	0.000	0.000	0.000	0.000	
co2p63	0.000	0.000	0.000	0.000	0.000	
co2p64	0.000	0.000	0.000	0.000	0.000	
co2p65	0.000	0.000	0.000	0.000	0.000	
co2p66	0.000	0.000	0.000	0.000	0.000	
co2p67	0.000	0.000	0.000	0.000	0.000	
co2p68	0.000	0.000	0.000	0.000	0.000	
co2p69	0.000	0.000	0.000	0.000	0.000	
co2p70	0.000	0.000	0.000	0.000	0.000	
co2p71	0.000	0.000	0.000	0.000	0.000	
co2p72	0.000	0.000	0.000	0.000	0.000	
co2p73	0.000	0.000	0.000	0.000	0.000	
co2p74	0.000	0.000	0.000	0.000	0.000	
co2p75	0.000	0.000	0.000	0.000	0.000	
co2p76	0.000	0.000	0.000	0.000	0.000	
co2p77	0.000	0.000	0.000	0.000	0.000	
co2p78	0.000	0.000	0.000	0.000	0.000	
co2p79	0.000	0.000	0.000	0.000	0.000	
co2p80	0.000	0.000	0.000	0.000	0.000	
co2p81	0.000	0.000	0.000	0.000	0.000	
co2p82	0.000	0.000	0.000	0.000	0.000	
co2p83	0.000	0.000	0.000	0.000	0.000	
co2p84	0.000	0.000	0.000	0.000	0.000	
co2p85	0.000	0.000	0.000	0.000	0.000	
co2p86	0.000	0.000	0.000	0.000	0.000	
co2p87	0.000	0.000	0.000	0.000	0.000	
co2p88	0.000	0.000	0.000	0.000	0.000	
co2p89	0.000	0.000	0.000	0.000	0.000	
co2p90	0.000	0.000	0.000	0.000	0.000	
co2p91	0.000	0.000	0.000	0.000	0.000	
co2p92	0.000	0.000	0.000	0.000	0.000	
co2p93	0.000	0.000	0.000	0.000	0.000	
co2p94	0.000	0.000	0.000	0.000	0.000	
co2p95	0.000	0.000	0.000	0.000	0.000	
co2p96	0.000	0.000	0.000	0.000	0.000	
co2p97	0.000	0.000	0.000	0.000	0.000	
co2p98	0.000	0.000	0.000	0.000	0.000	
co2p99	0.000	0.000	0.000	0.000	0.000	
co2p100	0.000	0.000	0.000	0.000	0.000	
co2p101	0.000	0.000	0.000	0.000	0.000	
co2p102	0.000	0.000	0.000	0.000	0.000	
co2p103	0.000	0.000	0.000	0.000	0.000	
co2p104	0.000	0.000	0.000	0.000	0.000	
co2p105	0.000	0.000	0.000	0.000	0.000	
co2p106	0.000	0.000	0.000	0.000	0.000	
co2p107	0.000	0.000	0.000	0.000	0.000	
co2p108	0.000	0.000	0.000	0.000	0.000	
co2p109	0.000	0.000	0.000	0.000	0.000	
co2p110	0.000	0.000	0.000	0.000	0.000	
co2p111	0.000	0.000	0.000	0.000	0.000	
co2p112	0.000	0.000	0.000	0.000	0.000	
co2p113	0.000	0.000	0.000	0.000	0.000	
co2p114	0.000	0.000	0.000	0.000	0.000	
co2p115	0.000	0.000	0.000	0.000	0.000	
co2p116	0.000	0.000	0.000	0.000	0.000	
co2p117	0.000	0.000	0.000	0.000	0.000	
co2p118	0.000	0.000	0.000	0.000	0.000	
co2p119	0.000	0.000	0.000	0.000	0.000	
co2p120	0.000	0.000	0.000	0.000	0.000	
co2p121	0.000	0.000	0.000	0.000	0.000	
co2p122	0.000	0.000	0.000	0.000	0.000	
co2p123	0.000	0.000	0.000	0.000	0.000	
co2p124	0.000	0.000	0.000	0.000	0.000	
co2p125	0.000	0.000	0.000	0.000	0.000	
co2p126	0.000	0.000	0.000	0.000	0.000	
co2p127	0.000	0.000	0.000	0.000	0.000	
co2p128	0.000	0.000	0.000	0.000	0.000	
co2p129	0.000	0.000	0.000	0.000	0.000	
co2p130	0.000	0.000	0.000	0.000	0.000	
co2p131	0.000	0.000	0.000	0.000	0.000	
co2p132	0.000	0.000	0.000	0.000	0.000	
co2p133	0.000	0.000	0.000	0.000	0.000	
co2p134	0.000	0.000	0.000	0.000	0.000	
co2p135	0.000	0.000	0.000	0.000	0.000	
co2p136	0.000	0.000	0.000	0.000	0.000	
co2p137	0.000	0.000	0.000	0.000	0.000	
co2p138	0.000	0.000	0.000	0.000	0.000	
co2p139	0.000	0.000	0.000	0.000	0.000	
co2p140	0.000	0.000	0.000	0.000	0.000	
co2p141	0.000	0.000	0.000	0.000	0.000	
co2p142	0.000	0.000	0.000	0.000	0.000	
co2p143	0.000	0.000	0.000	0.000	0.000	
co2p144	0.000	0.000	0.000	0.000	0.000	
co2p145	0.000	0.000	0.000	0.000	0.000	
co2p146	0.000	0.000	0.000	0.000	0.000	
co2p147	0.000	0.000	0.000	0.000	0.000	
co2p148	0.000	0.000	0.000	0.000	0.000	
co2p149	0.000	0.000	0.000	0.000	0.000	
co2p150	0.000	0.000	0.000	0.000	0.000	
co2p151	0.000	0.000	0.000	0.000	0.000	
co2p152	0.000	0.000	0.000	0.000	0.000	
co2p153	0.000	0.000	0.000	0.000	0.000	
co2p154	0.000	0.000	0.000	0.000	0.000	
co2p155	0.000	0.000	0.000	0.000	0.000	
co2p156	0.000	0.000	0.000	0.000	0.000	
co2p157	0.000	0.000	0.000	0.000	0.000	
co2p158	0.000	0.000	0.000	0.000	0.000	
co2p159	0.000	0.000	0.000	0.000	0.000	
co2p160	0.000	0.000	0.000	0.000	0.000	
co2p161	0.000	0.000	0.000	0.000	0.000	
co2p162						

CO₂ per kWh: Yeo-Johnson Transformed Distribution

Creation of a target variable for CO₂ per unit of energy (CO₂/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_normed['CO2_per_kWh'] = df_normed['y1_CO2'] / df_normed['Usage_kWh']  
#Exploring the distribution. Must be aware the outliers and skew and other stuff that might be going on  
import seaborn as sns  
plt.figure(figsize(8, 4))  
sns.histplot(df_normed['CO2_per_kWh'], bins=30, kde=True, color="pinkred")  
plt.title("Distribution of CO2 per kWh (Yeo-Johnson Transformed)")  
plt.xlabel("y1_CO2 / y1_kWh")  
plt.tight_layout()  
plt.show()
```



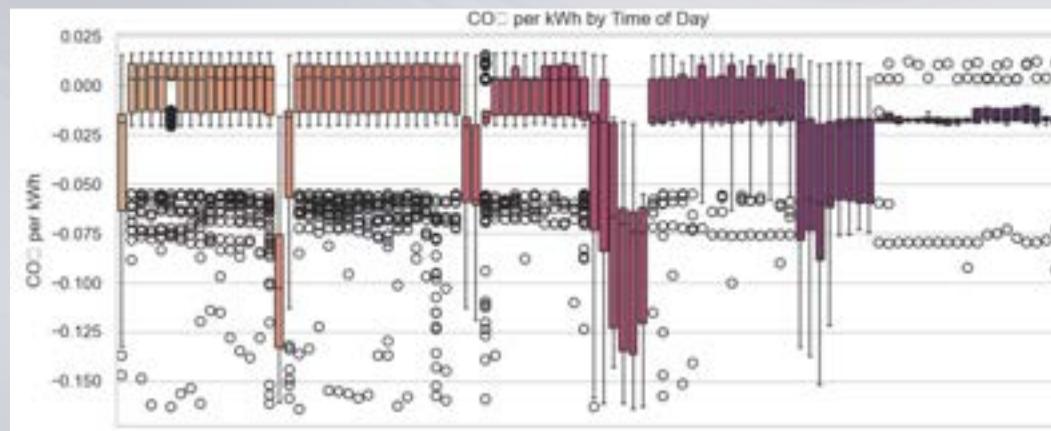
The transformed distribution revealed several distinct behaviors:

- 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 CO₂_per_kWh by Time of Day



This boxplot reveals how emissions intensity (CO₂/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.

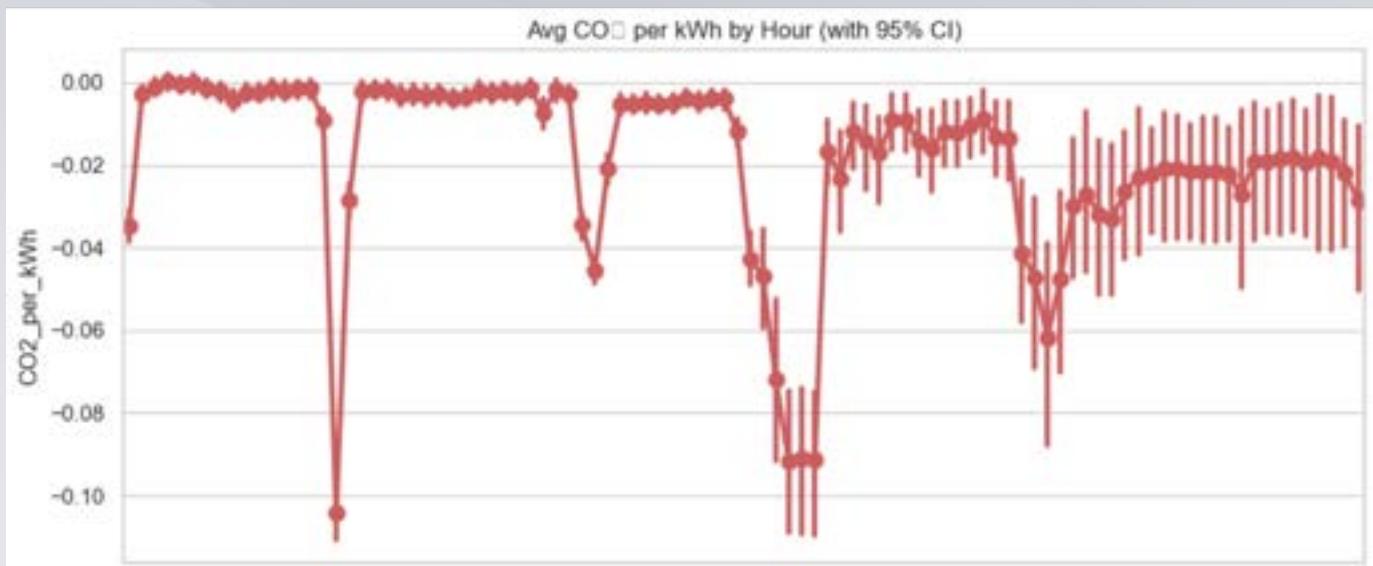


Observations from Plot

- 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 CO₂ per kWh with Confidence Bands

To validate the time-dependent behavior observed in the boxplots, I generated a point plot of average CO₂ 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

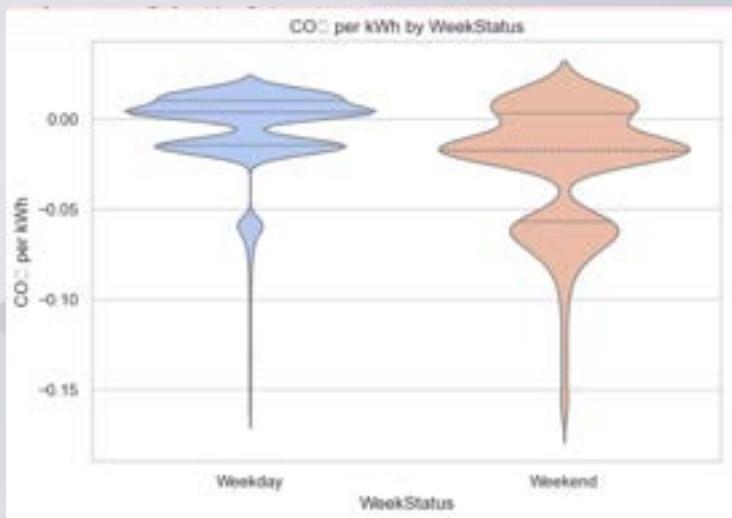


Violin Plot of CO₂ per kWh Distribution by Day Type

Exploring Weekday vs. Weekend Emission Profiles

Given the temporal sensitivity observed in both the CO₂ 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

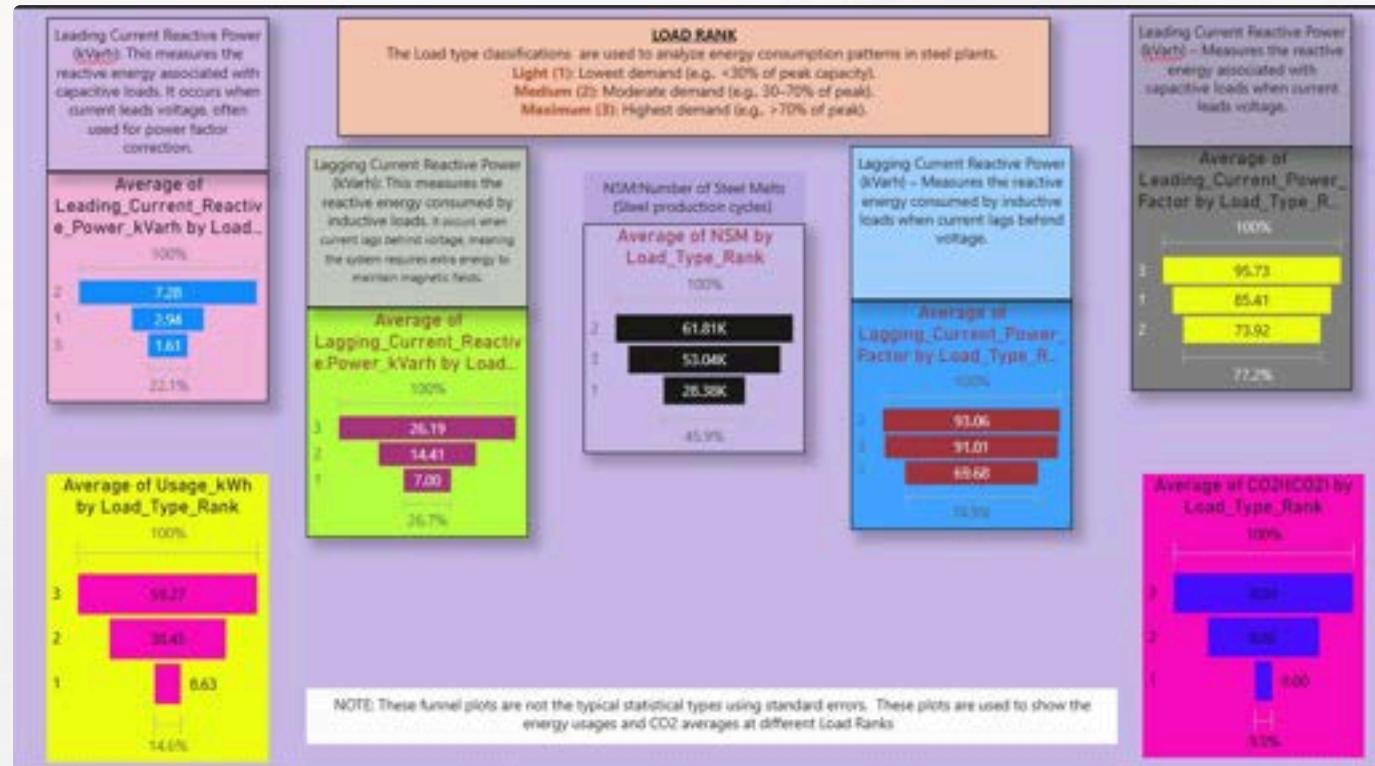


Weekdays show a wider top density near zero, narrowing quickly with a brief blip around -0.05, then tapering into a thin tail toward -0.15. This suggests a central concentration of operating states, with tightly regulated emissions and fewer extreme outliers.

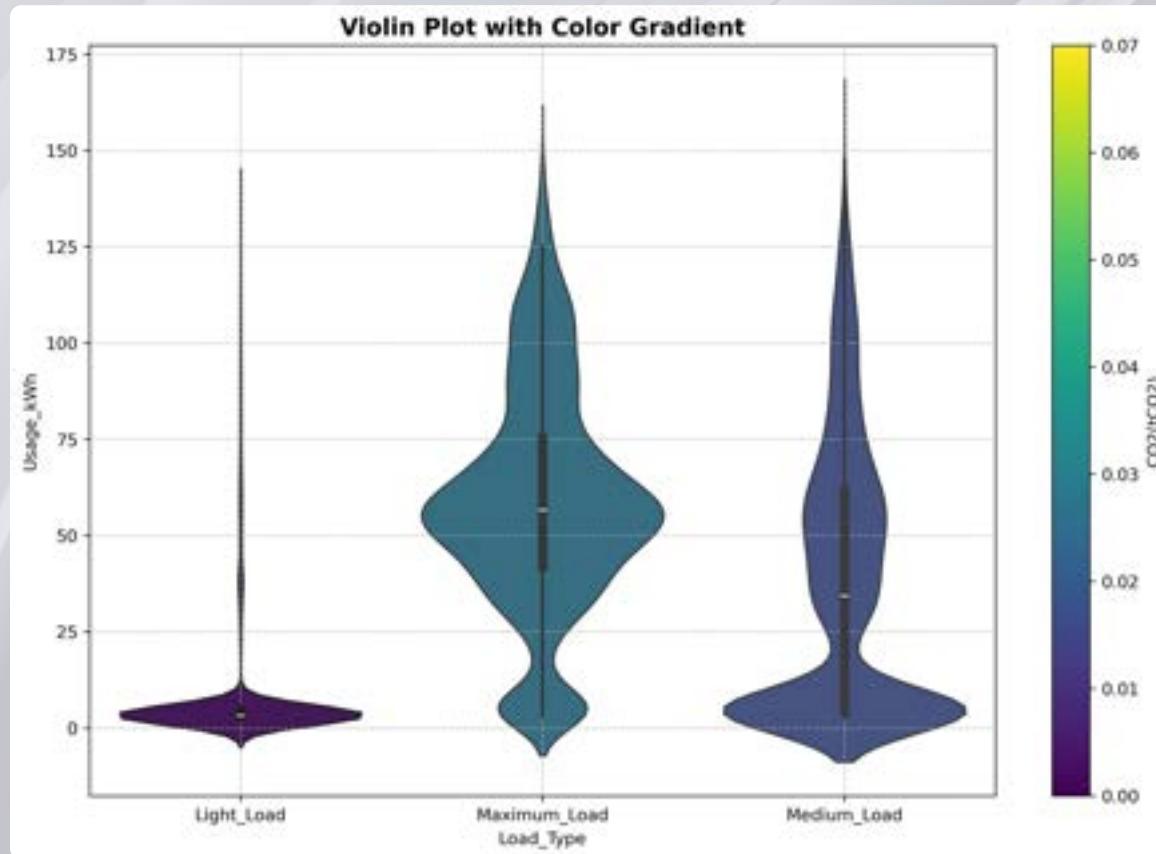
Weekends, in contrast, feature a slightly flatter top, followed by a broader second segment, especially near -0.05, where density thickens. Rather than tightening into a fine tail like weekdays, the weekend distribution remains more dispersed and rounded, implying diverse operational conditions or less 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 CO₂ 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.

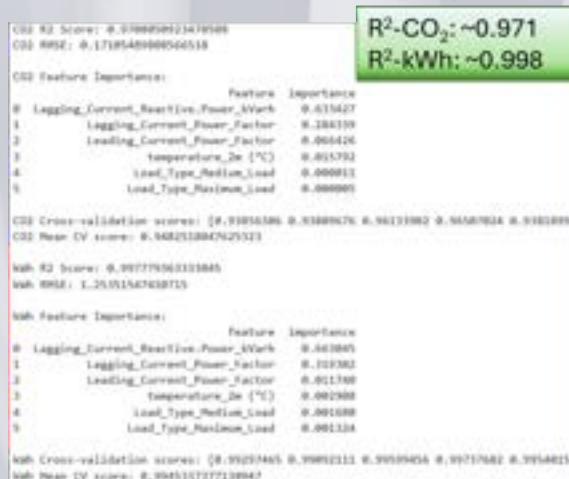
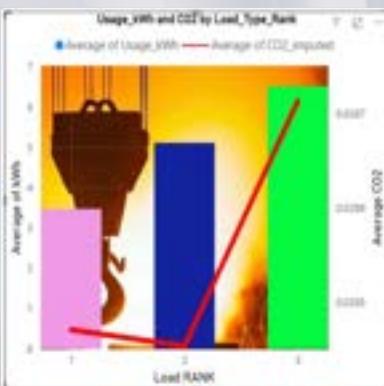
These distributional cues preview what the model will later confirm: operational load rank isn't just a label; it's a behavioral fingerprint.

Random Forest Model with Categorical Load



Load Rank Revisited: Operational Demand Meets Emission Behavior

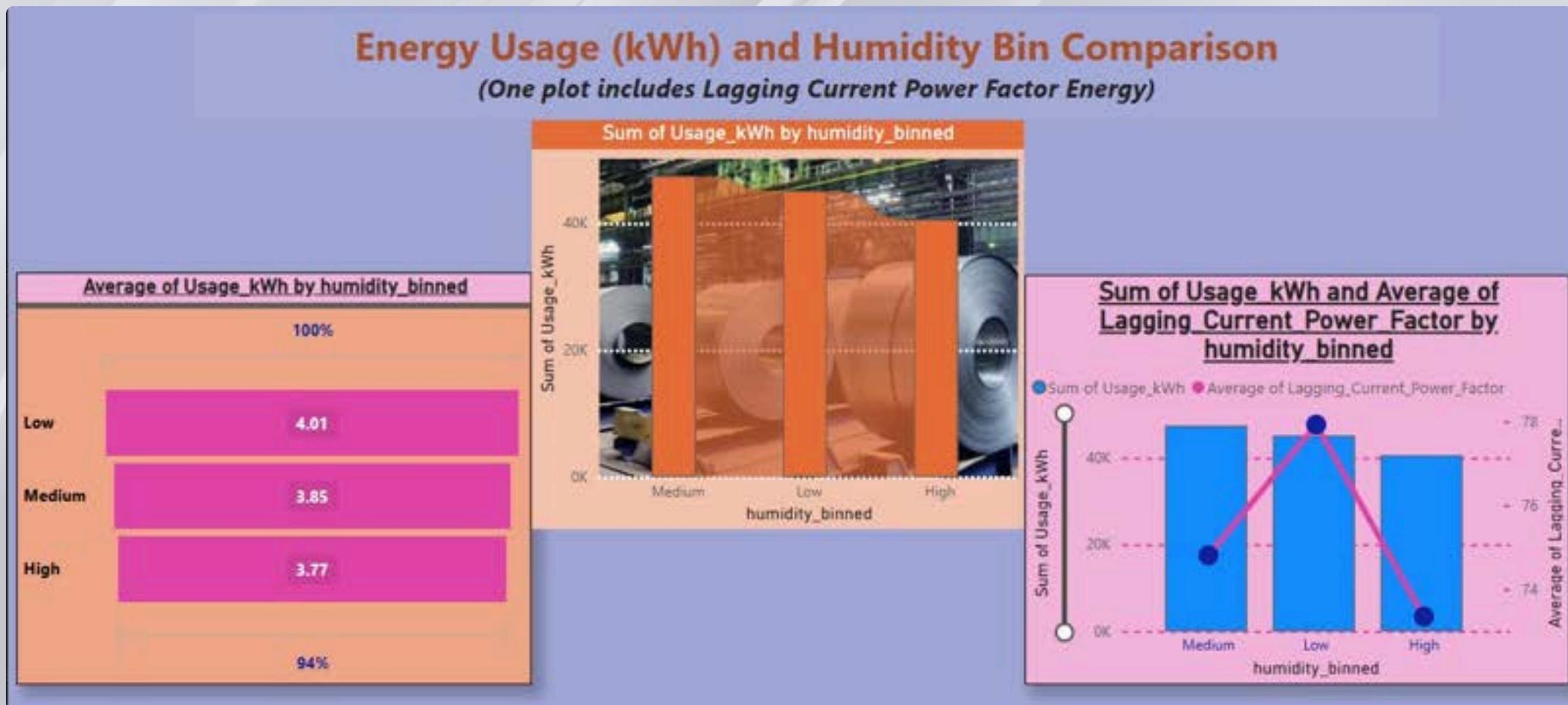
Earlier in my analysis, Load Rank emerged as a key operational marker with distinguishing low, medium, and peak production states. I examined its relationship to energy consumption and emissions behavior, uncovering nonlinear shifts and efficiency transitions. Now, using Random Forest models tailored to both kWh and CO₂, 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.



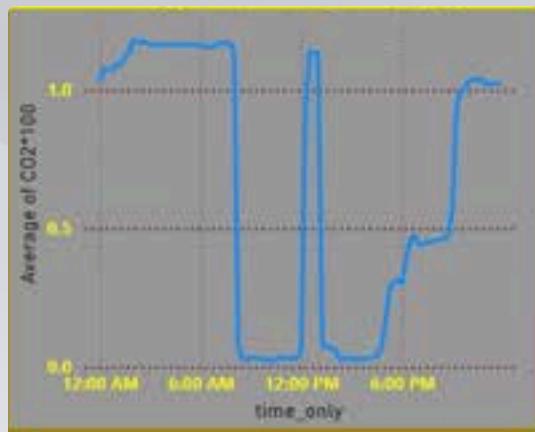
Interpretive Highlights

- **Reactive metrics dominate** in both models; the system's internal load dynamics are the clear drivers.
- **Temperature matters modestly**, more so for CO₂ 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.

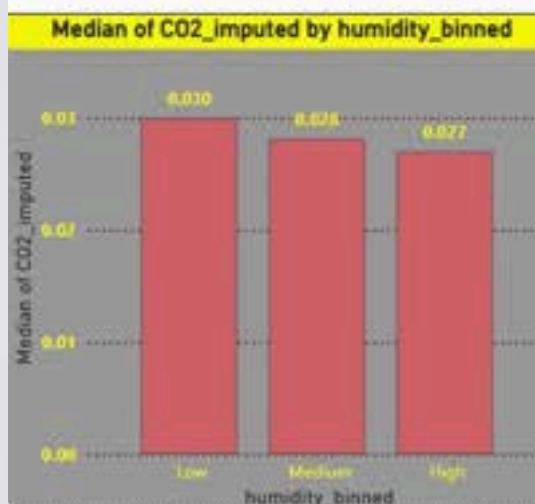
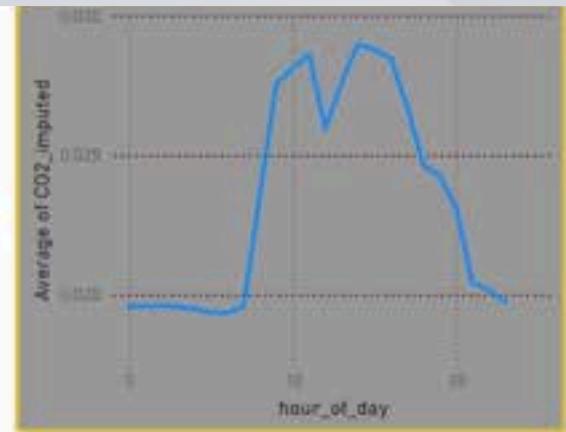
While categorical load patterns revealed operational structure through the RF model, system behavior doesn't unfold in isolation. Environmental conditions, particularly humidity, introduce subtle but measurable shifts in energy usage and emissions. Zooming into these climate driven regimes, it was uncovered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



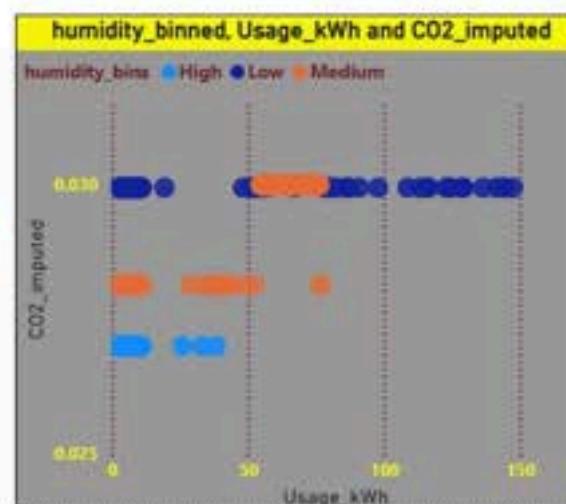
Comparison of CO₂ Comparison of actual and imputed CO₂ emissions across time, humidity levels, and energy usage., Imputed CO₂ against Time and Humidity and Energy.



Imputing to take care of missing data



Adding Usage kWh, (used in models)



HUMIDITY: Diagnostically justified, not just Exploratory

Dunn Test: CO₂ and Humidity Bins

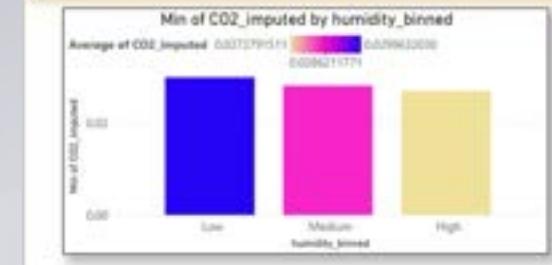
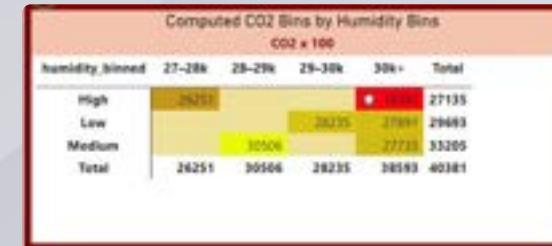
With p-values so tiny the takeaway is loud and clear: each humidity bin's CO₂ distribution is significantly different from the others. Not just statistically significant, but celestially significant if that existed. This validates what a Kruskal-Wallis test hinted at earlier: The humidity's effect isn't subtle. It appears that the imputed CO₂ 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 CO₂.

```
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	



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)	Directionality & baseline insights	Great for quick diagnostics, identifying linear relationships
VECM	Cointegration-based system modeling	Temporal equilibrium & system feedback	Ideal for long-run relationships and shock-response analysis
Random Forest	Ensemble tree voting (bagging)	Robust prediction & feature interpretability	Excellent with SHAP, ICE, PDPs for uncovering nonlinear drivers
XGBRegressor	Gradient-boosted decision trees (iterative error correction)	High-performance modeling	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 (v0.4) - Keeping track and saving ...
from sklearn import Ridge
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\srinivas\OneDrive\Desktop\train_low_humidity_cleaned.csv')
print("Dataset shape:", df.shape)
print("Low Humidity rows?", len(df[df['humidity_cleaned'] == 'low'])))

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

# Define features (Include log_kWh, exclude CO2_lag_1 to avoid usage_MW)
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_cleaned'] == '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 (unedited parameters)
model_low = Ridge()
model_low.fit(X_train, y_train)

# Predict
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))

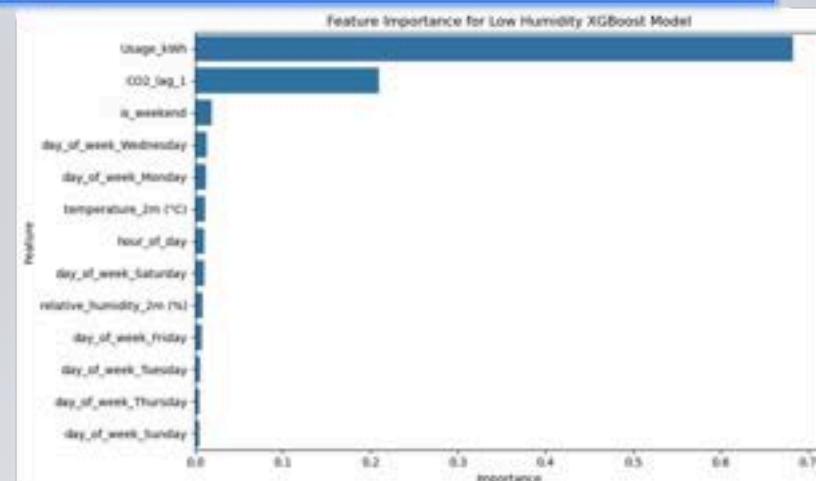
# Save feature importance
importances = pd.DataFrame(['Feature': features, 'Importance': model_low.feature_importances_])
print('Importances')
print(importances.sort_values('Importance', ascending=False))
importances.to_csv(r'C:\Users\srinivas\OneDrive\Desktop\feature_importance_low_gbtree_v0.4.csv', index=False)
print('Feature importance saved to feature_importance_low_gbtree_v0.4.csv')
```

Overview and Thoughts

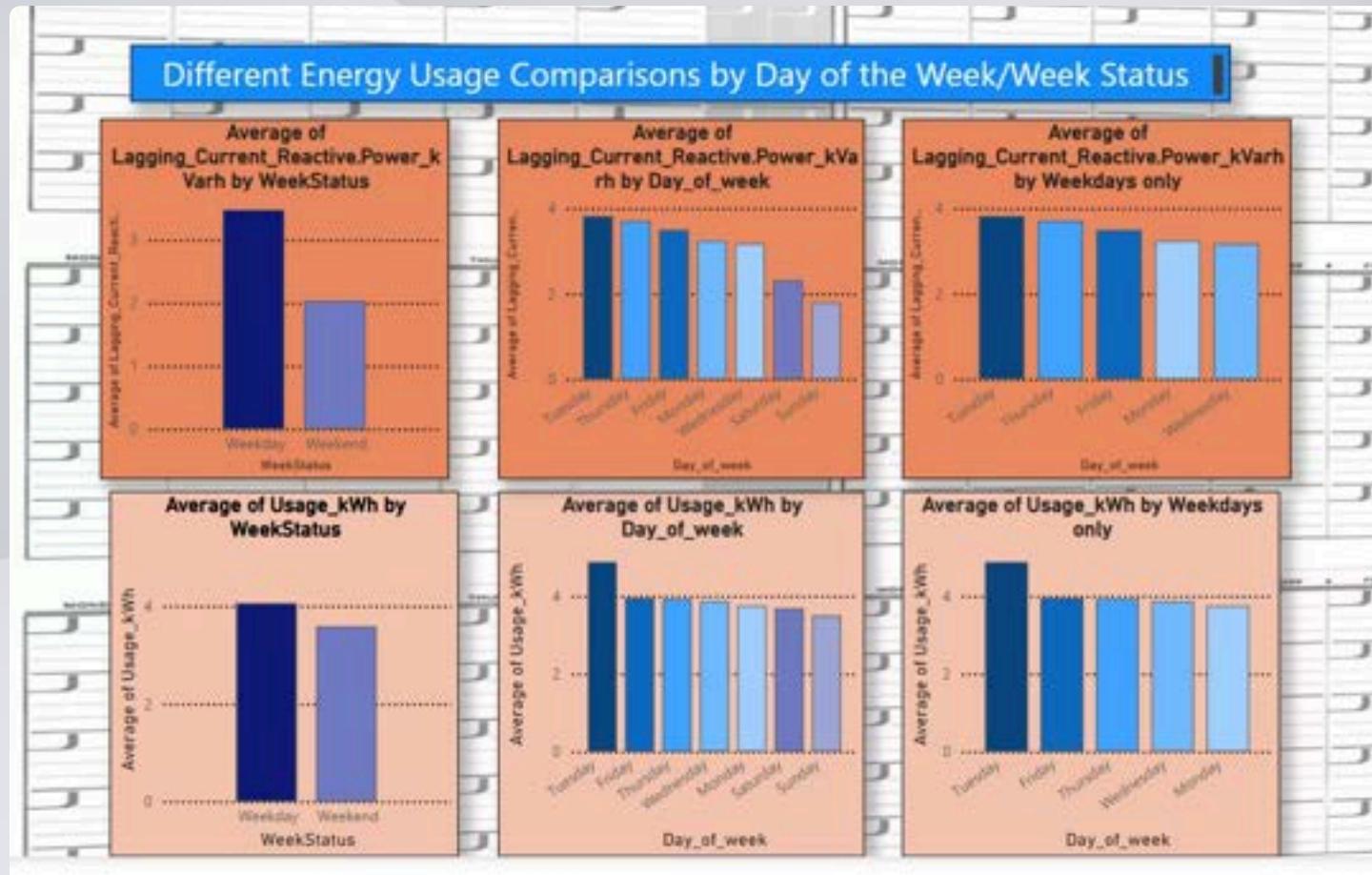
- **Excellent predictive strength:** $R^2 \sim 0.989$ on both train and test sets with microscopic MSE
- **Log transformation stabilizes variance and normalizes skew:** Given earlier diagnostics, this aligns beautifully with the system's distributional trait which is especially helpful under low humidity 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.

Feature variances:

log_kWh	1.420569
temperature_2m (°C)	117.507486
relative_humidity_2m (%)	122.232102
hour_of_day	34.667253
is_weekend	0.217252
day_of_week_Friday	0.141803
day_of_week_Monday	0.115526
day_of_week_Saturday	0.140586
day_of_week_Sunday	0.127354
day_of_week_Thursday	0.127430
day_of_week_Tuesday	0.097780
day_of_week_Wednesday	0.103319



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

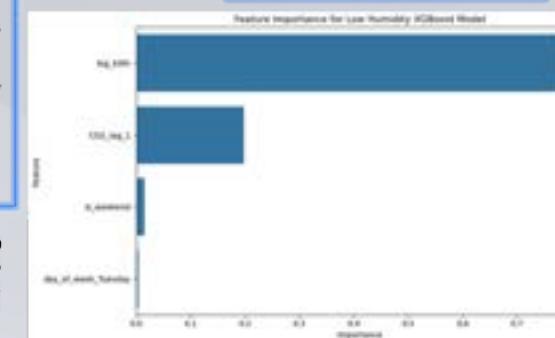
XGBoost Forest Model with log_kWh With limited Temporal Variables and Lag

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.

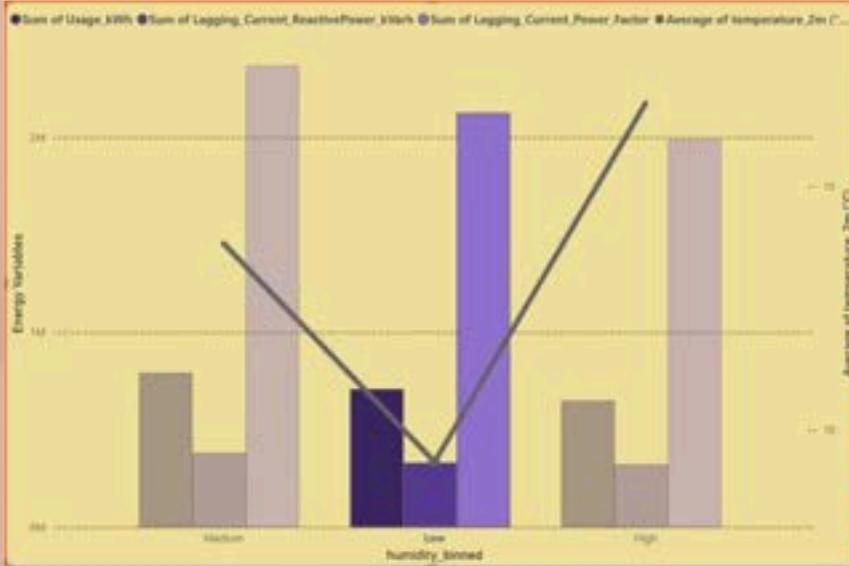
Low Humidity Train R²: 0.9865466799168923
Low Humidity Test R²: 0.990629537339741
Low Humidity Train MSE: 2.0169477900957362e-0
Low Humidity Test MSE: 1.3240912660866145e-07

Overview and Thought

- This model hits nearly perfect R² with just four features which is proof the system's behavior 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..



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.

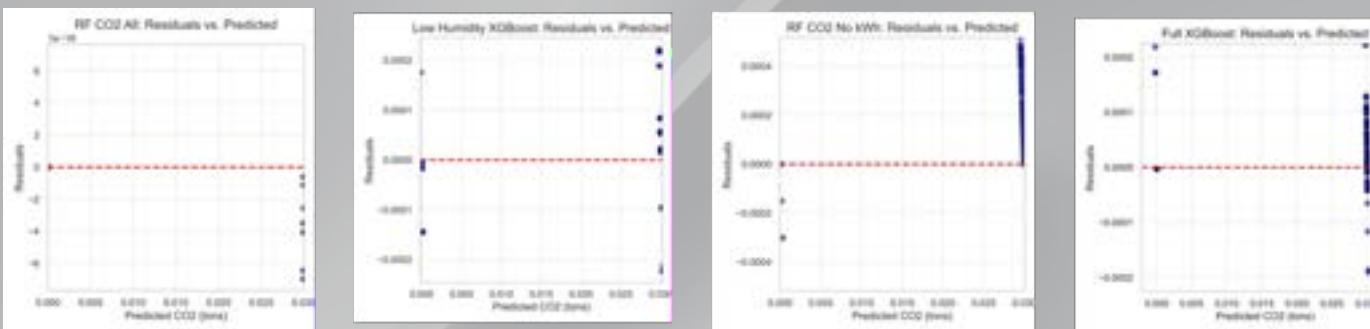
Unified Variable Suite — Model Interpretation Side-by-Side



Model Performance Comparison : CO₂ Prediction

Model	R ² Score	RMSE	Notes
RF (All Variables)	0.9985	0.0004362	Tight error spread, minimal outliers — robust baseline
XGBoost (Low Humidity)	0.9944	0.0009868	Best Q-Q fit, fewer outliers — strong under dry regime
Full XGBoost	0.9951	0.0007801	Slight right-tail curvature, moderate outlier density
RF (No Usage_kWh)	0.9944	0.0008356	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.

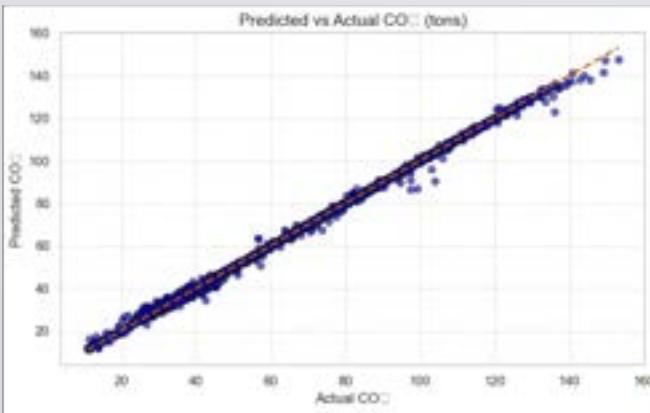


Plot 1: RF (All Variables) shows symmetrical error with two edge outliers, indicating strong predictive balance. **Plot 2: XGBoost (Low Humidity)** exhibits the best distribution fit with minimal deviation from the normal curve. **Plot 3: RF (No Usage_kWh)** reveals a wider error spread and structural 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

```
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()
```



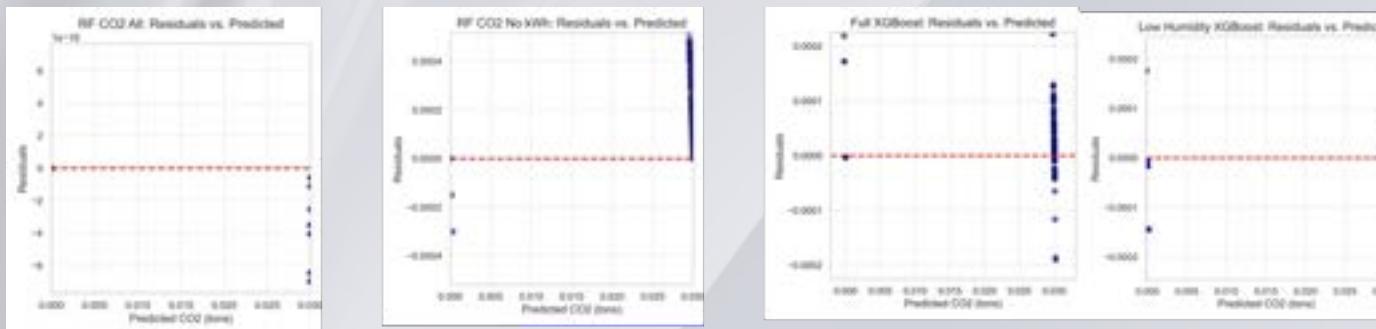
While nonlinear models captured complexity across regimes and environmental layers, the final OLS regression delivered elegant simplicity. With well-behaved inputs and a refined feature set, the predicted vs. actual CO₂ plot shows near-perfect alignment with each point hugging the diagonal line. This fit reflects not just statistical performance, but how far the modeling journey has come—from messy signals to interpretable clarity.

Residual Plot Interpretation: CO₂ 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

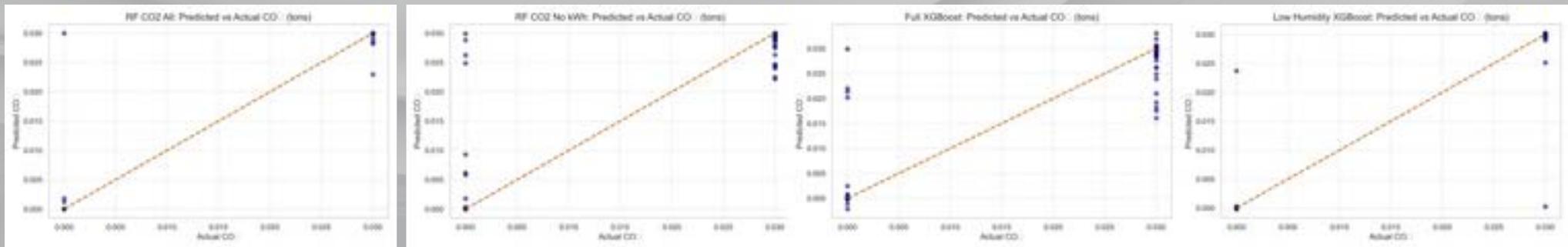
Each residual plot doesn't just reveal error; it reveals model personality.

- 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

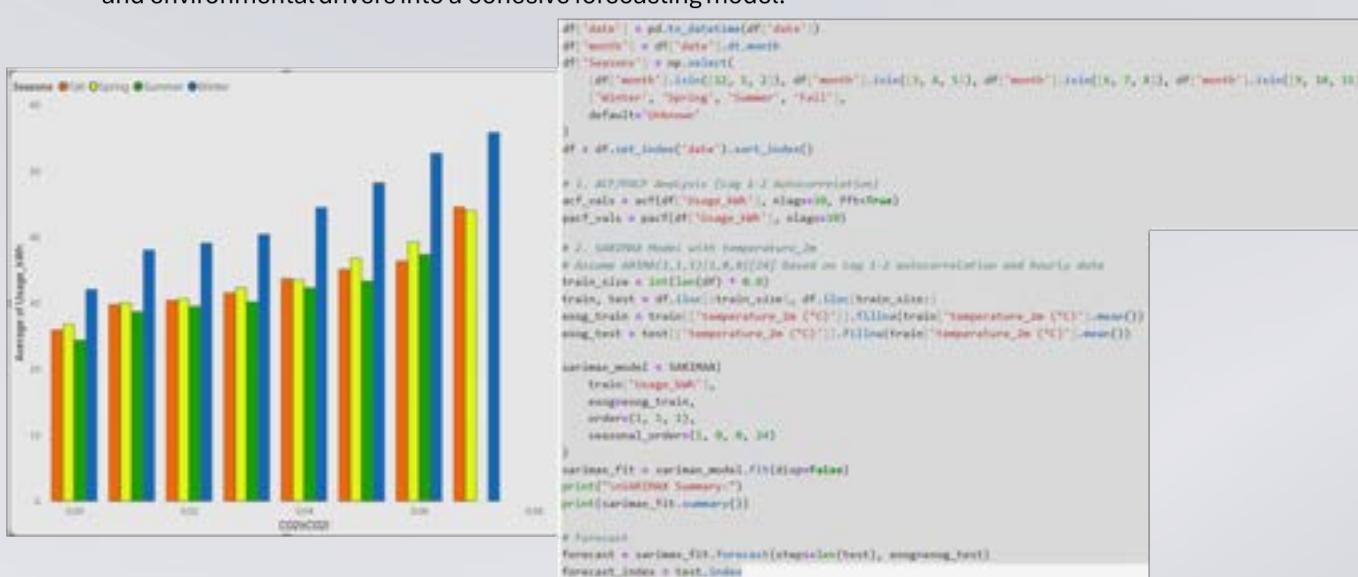
Model	Residuals Insight	Q-Q Plot Insight
RF (All Variables)	Tightest residual spread, single outlier on each end	Slight curvature at tails with closest fit to normal distribution
Low Humidity XGB	Slightly more residual variance, but no curvature	Nearly perfect Q-Q alignment, minimal outliers and best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



SARIMAX and Seasonality Results

Overview and Thoughts

- **Strong temporal structure:** The AR and MA components confirm energy usage is highly autocorrelated and responds to recent history.

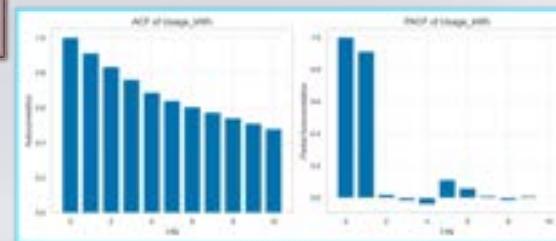
- **Seasonality is real:** Winter months show elevated demand and Kruskal-Wallis confirms this statistically. This supports previous analysis regarding winter months and energy usage.

- **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

		Latent Needs			
Item	Variable	Intercept	Slope	Mean	SD
Latent Needs	Latent_Needs_1, L_1, L_2,L_X, L_M	-0.00000000	0.00000000	0.00000000	0.00000000
Time	Time, All but 1980	-0.00000000	0.00000000	0.00000000	0.00000000
Sample	01-01-1980-4021	-0.00000000	0.00000000	0.00000000	0.00000000
Time	19-19-1980	-0.00000000	0.00000000	0.00000000	0.00000000
Measurement Type	NA				
	coeff	std. err.	t	F(1,1)	(N, R)
Intercept, M	2.0000	0.120	16.668	0.0000	1.000
pr_12	0.2000	0.060	3.333	0.0000	0.000
pr_13	0.2000	0.060	3.333	0.0000	0.000
pr_14	0.2000	0.060	3.333	0.0000	0.000
signed	100.0000	0.750	133.333	0.0000	107.100
standard error					
Using obs 1-111 (1)	14.00	Surface Area (cm²)		0.000000.10	
Model 1	0.00	Prob>F(1)		0.00	
Residual standard error:	0.00	Mean		0.00	
Number of observations:	0.00	Standard Deviation		0.00	
LRTDRX: Residual Statistics:					
count	1000000000				
mean	0.00000000				
sd	0.00000000				
se	0.00000000				
skew	-0.00000000				
kurt	-0.00000000				
min	-12.00000000				
max	12.00000000				
range	24.00000000				
df	1				
Pr>ChiSq	0.00000000				
Frank-Wallis Test: Bartlett's χ^2 = 330.13, p-value = 0.000000					

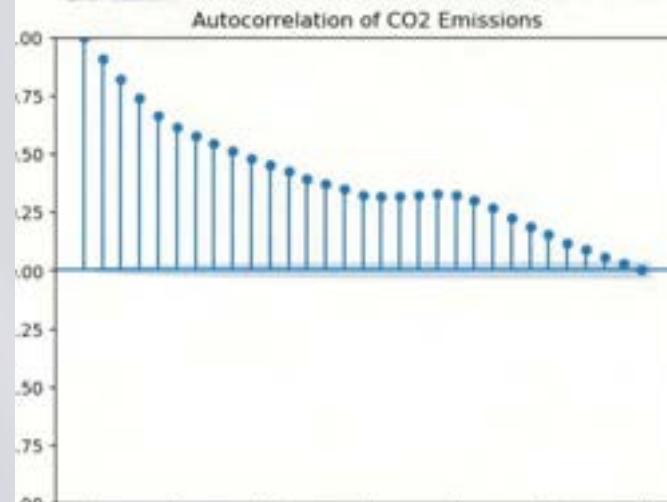
Metric / Parameter	Value / Description
Dependent Variable	Usage_kWh
Observations	28,032 hourly records (Jan–Oct 2018)
SARIMAX Order	ARIMA (1,1,1) × (1,0,0)[24]
Exogenous Variable	temperature_2m (°C)
Log Likelihood	-113,950.60
AIC / BIC / HQIC	227,911.20 / 227,952.41 / 227,924.46
AR(1) Coefficient	0.8961 (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



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

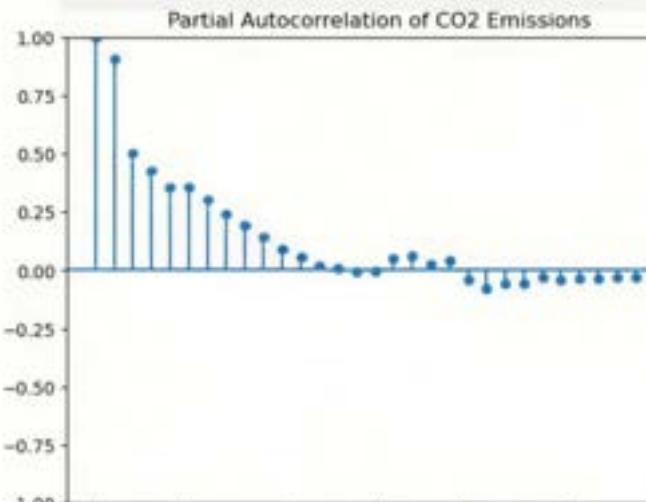
# Load emissions data

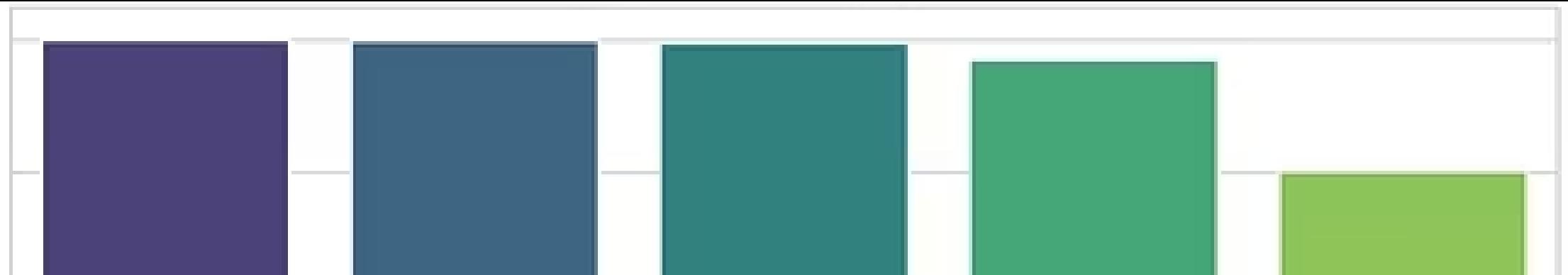
plot_acf(df["CO2(tCO2)"], lags=30) # Adjust lags based on data
plt.title("Autocorrelation of CO2 Emissions")
plt.show()
```



```
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()
```





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

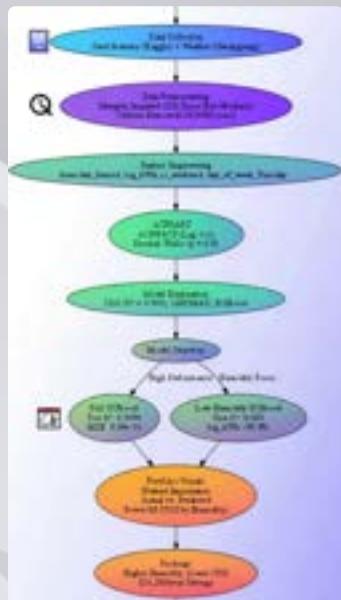
log_kWh contributed 98.9% to prediction
accuracy

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
OLS	Linear Regression (Ordinary Least Squares)	Interpretable Relationships; Directionality & Baseline Insights	Great for quick diagnostics, identifying linear relationships	Final fit showed mixed fit predictions and excellent baseline clarity
RF	Ensemble trees	Nonlinear structure, Robust Prediction & Feature importance/Interpretability	Using SHAP, ICE, PDPs for unpacking nonlinear drivers is insightful	Strong residual diagnostics; categorical load modeling with insight
XGBoost	Gradient-boosted trees	High-Performance Modeling Under Regressions	Captures subtle interactions, can excel in tuned prediction tasks	Bent Q-Q fit under low humidity; sensitive to weekly patterns
VECM	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 relevance
SARIMAX	Seasonal time series with exogenous factors	Short-term forecasting with rhythm, time 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 ₂ Prediction			
Model	R ² 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 ² 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 CO₂ 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 CO₂. 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 CO₂ & 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 CO₂ 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 CO₂ 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.

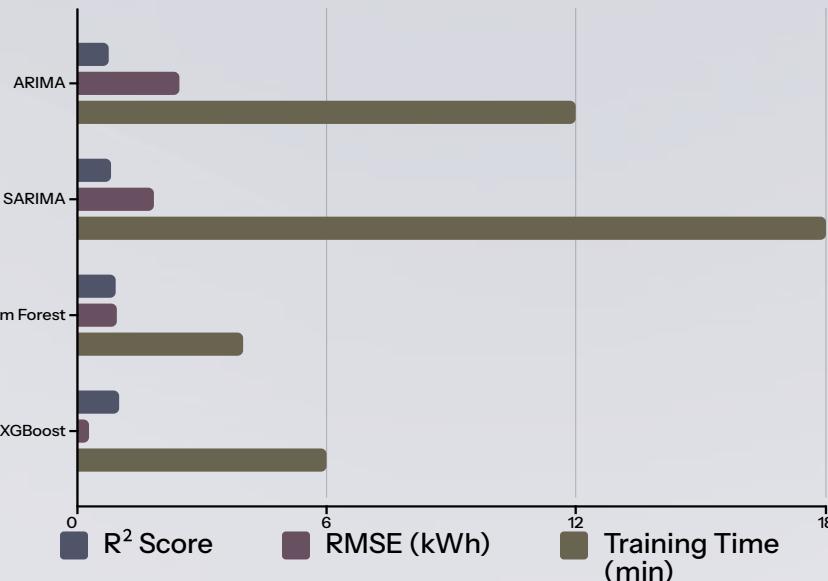
These structural models laid the foundation for temporal analysis. Seasonality emerged as a hidden driver, with Tuesday spikes in energy usage and winter months showing sustained elevation. To capture this rhythm, I pivoted to SARIMAX modeling, allowing me to fold in monthly effects, autocorrelation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

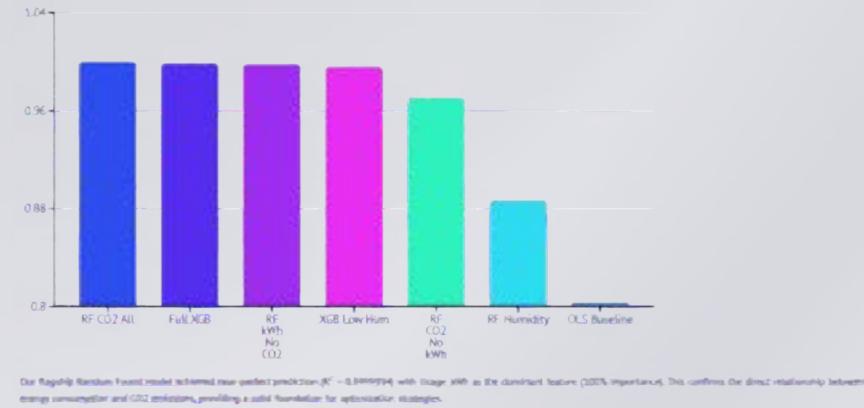
Model Performance Analysis

Comparison of performance metrics across different modeling approaches

Performance Metrics Comparison



Model Performance Overview



XGBoost demonstrated superior performance with the highest R² 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.

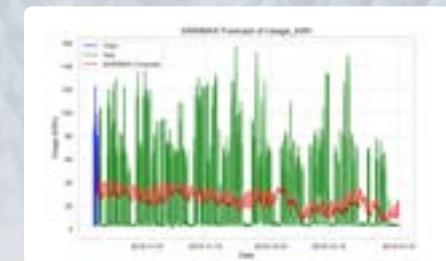
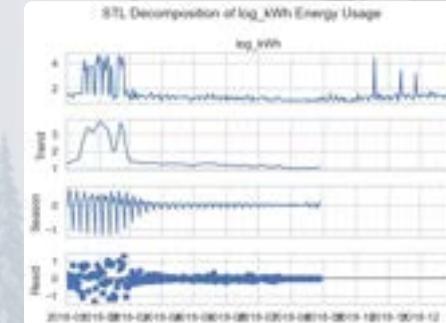
Model Approach: Decoding System Behavior Over Time

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



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

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.



Energy Through Time: A Dynamic System

Visualizing the steel plant's energy consumption and CO₂ 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

CO₂, kWh, and Temporal Forecasting

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

¹ RMSE in log scale ² Estimated from earlier decomposition and fit metrics

Modeling with Structure, Signal & Seasonality

Interpretive Highlights

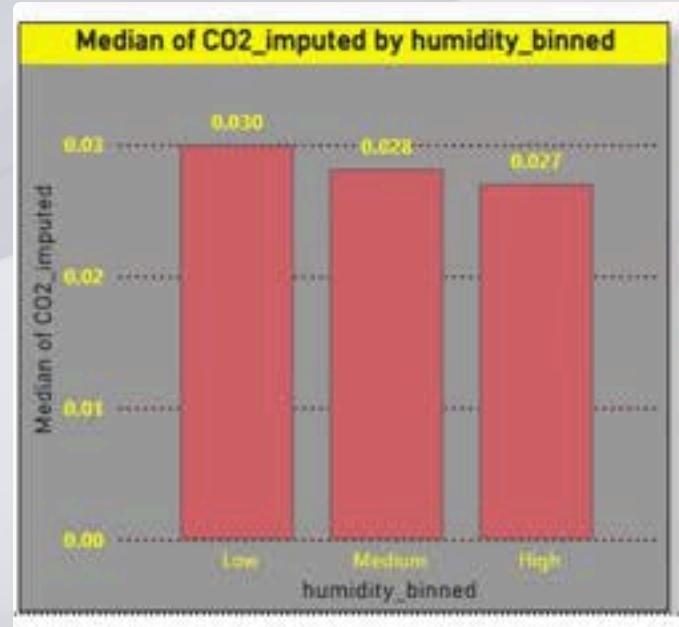
- Usage_kWh remains the system heartbeat, whether as a feature or forecast target.
- Low humidity reveals clarity that models tighten their signal when environmental noise recedes.
- Temporal models (SARIMAX, ARIMA) offer insight into cyclical demand, especially around winter peaks and weekday rhythm.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

Final Takeaways

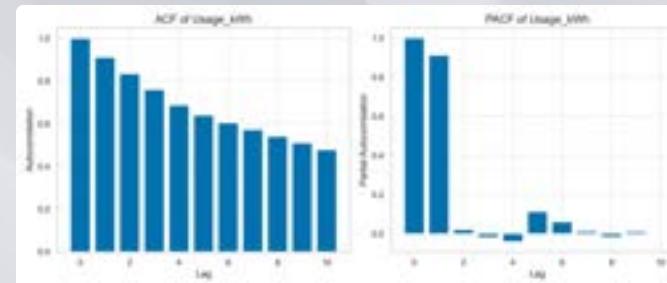
- Usage_kWh isn't just predictive, it's structural. It anchors both CO₂ and kWh models and reveals behavioral regimes.
- Temporal features (weekday, weekend, CO₂ lag) layers complexity and help low humidity models “sing.”
- SARIMAX introduces rhythm, seasonal intuition, and exogenous awareness. It's 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

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.

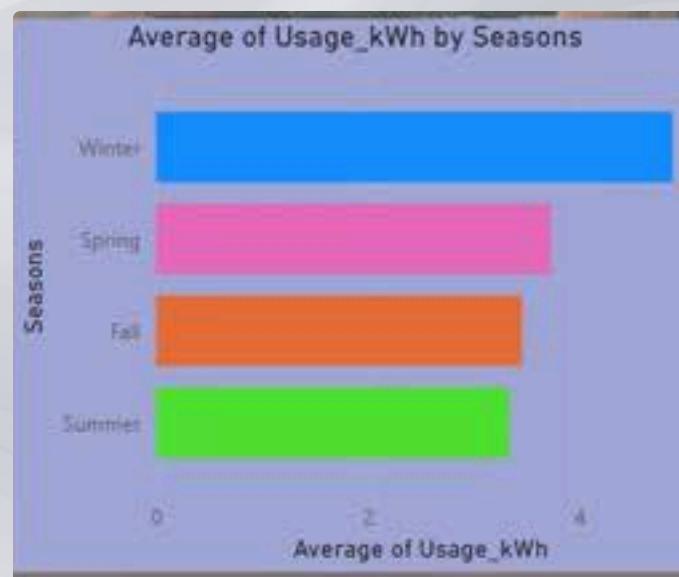
ACF Plot

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

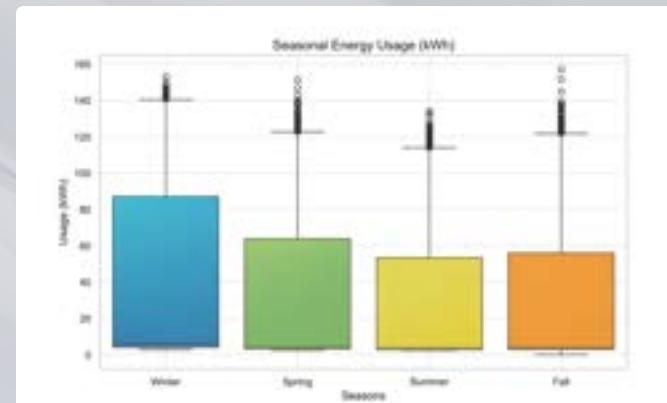
PACF Plot

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

Winter peak energy usage, driving savings focus.



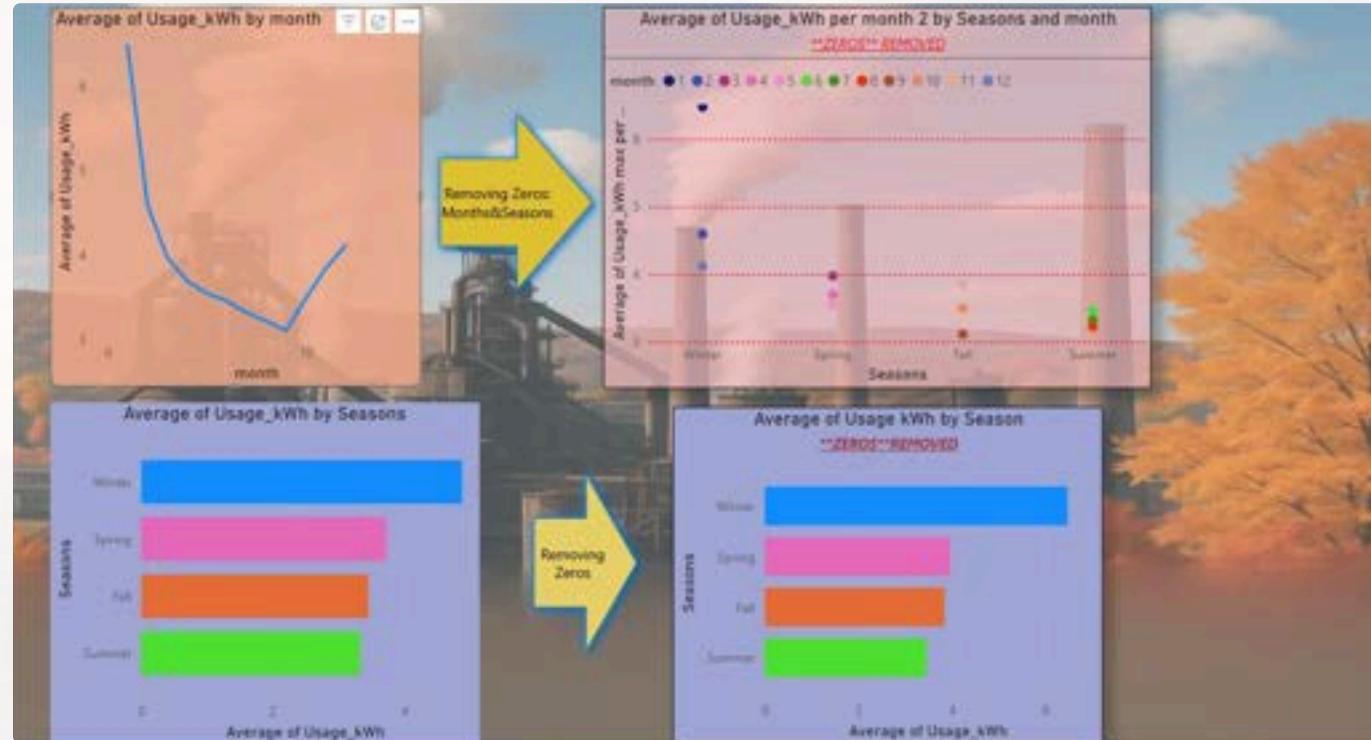
Kruskal-Wallis test confirmed statistically significant seasonal variation in energy usage ($p < 0.001$).**



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

tual vs. Predicted Energy Consumption

Energy Usage by Season: Winter Wins, unless looking at Costs!



Section 7: Potential Cost Savings

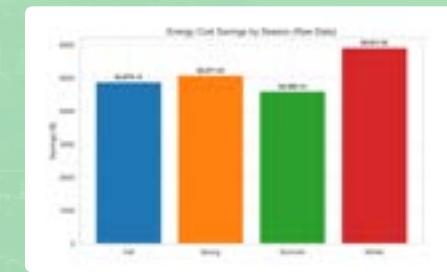




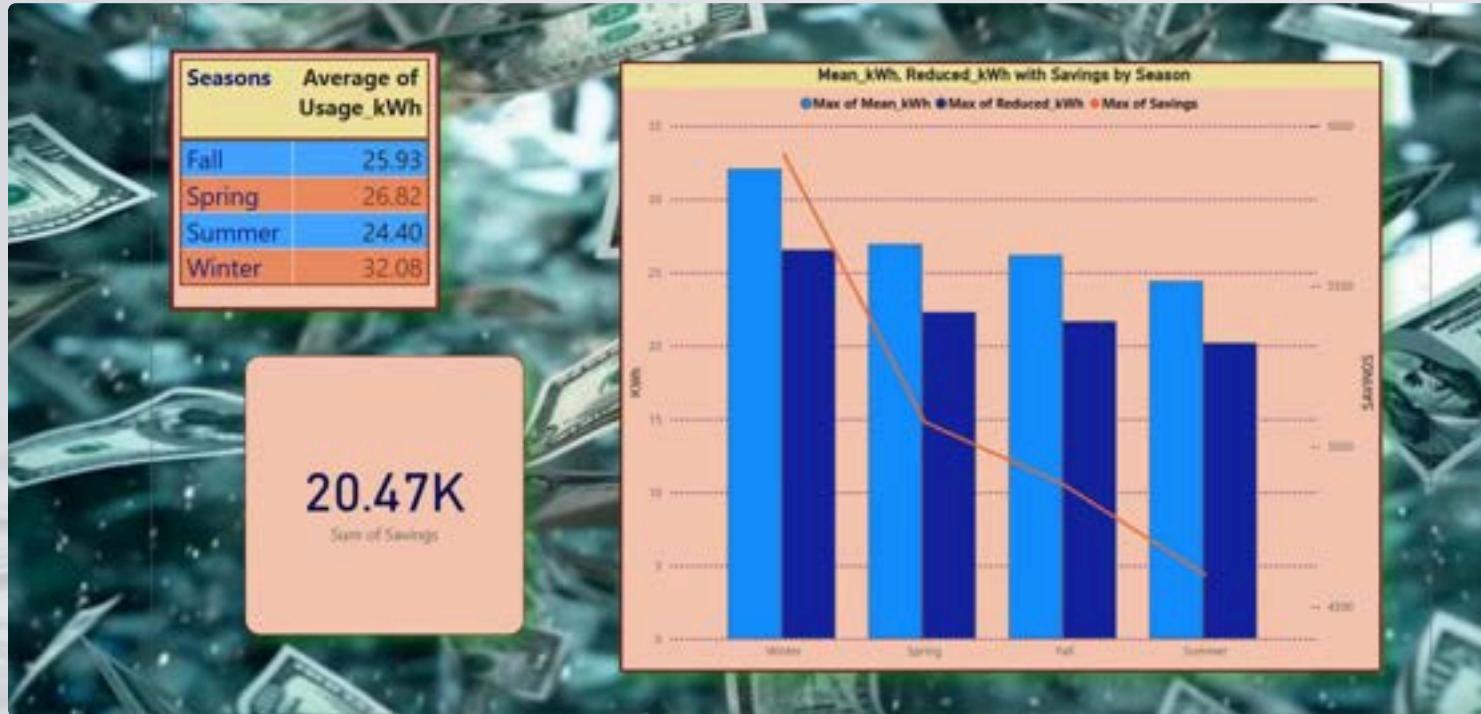
Sometimes, It's All About the Bottom Line

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

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



The Seasons and the Energy and Savings Comparison



CONCLUSIONS: Curves Were Fit, CO₂ 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 CO₂ 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

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

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