# Statistical Analysis of Retail Electricity Sales in the United States

*A Multiregional, Multisector Time Series Evaluation (2014–2024)*

## Mohankrishna Kolla

December 5, 2025

**Abstract**

**Abstract.** This report investigates temporal, seasonal, sectoral, and regional variation in retail electricity sales across the United States using archival monthly energy data from 2014–2024. Using exploratory data analysis, grouped descriptive statistics, and multiple linear regression on log-transformed sales, the study evaluates structural differences across regions and sectors. Results indicate strong statistical significance in regional and sectoral variation ($R^2 \approx 0.852$), with the South region and Residential sectors driving peak demand. The Transportation sector exhibits minimal grid impact relative to other sectors. Implications for demand forecasting and energy planning are discussed.

# Contents

# 1   Research Question and Motivation

The central research questions of this study are:

1. How do retail electricity sales differ across economic sectors (Residential, Commercial, Industrial, Transportation)?

2. How do these sectoral differences vary across U.S. Census regions (Northeast, Midwest, South, West)?

3. Is there evidence of a long-term national trend in electricity consumption over the 2014–2024 period?

4. How do seasonal patterns manifest in aggregate electricity demand?

Understanding these relationships is essential for long-term demand forecasting, infrastructure planning, and regional energy policy. Variation in electricity use reflects climate, demographics, regional industry, and transportation electrification.

# 2   Methods

## 2.1   Data Source and Preparation

This project uses publicly available archival monthly electricity sales data obtained from the U.S. Energy Information Administration (EIA). The dataset includes monthly observations of retail electricity sales (MWh) categorized by State and Sector.

Key preprocessing steps included:

- **Filtering:** Removal of aggregate "Total" rows to prevent double counting.

- **Feature Engineering:** Mapping states to Census Regions (Northeast, Midwest, South, West) and extracting seasonal indicators.

- **Transformation:** Application of a log-transformation ($\log(1 + \text{Value})$) to address heteroscedasticity and the heavy right-skew of electricity sales data.

- **Sampling:** A subset of 5,000 observations was used for the regression model to ensure computational efficiency while maintaining statistical power.

# 3 Exploratory Data Analysis

## 3.1 Overall Summary Statistics

Table 1 presents the summary statistics for the filtered dataset. The data exhibits high variance, with a standard deviation exceeding the mean in the original scale, necessitating the log transformation used in the inferential modeling.

**Table 1:** Summary statistics for retail electricity sales (Original vs. Log-Transformed).

| Variable | Count | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| Original Value (MWh) | 25,695 | 1,325.93 | 1,489.76 | 0 | 8,499.00 |
| Log-Transformed | 25,695 | 5.54 | 2.85 | 0 | 9.05 |

## 3.2 Overall Temporal Pattern

Figure 1 shows the mean electricity sales over time, averaged across all states, regions, and sectors. The data show pronounced seasonal patterns with consistent annual cycles and relatively stable long-run average levels.
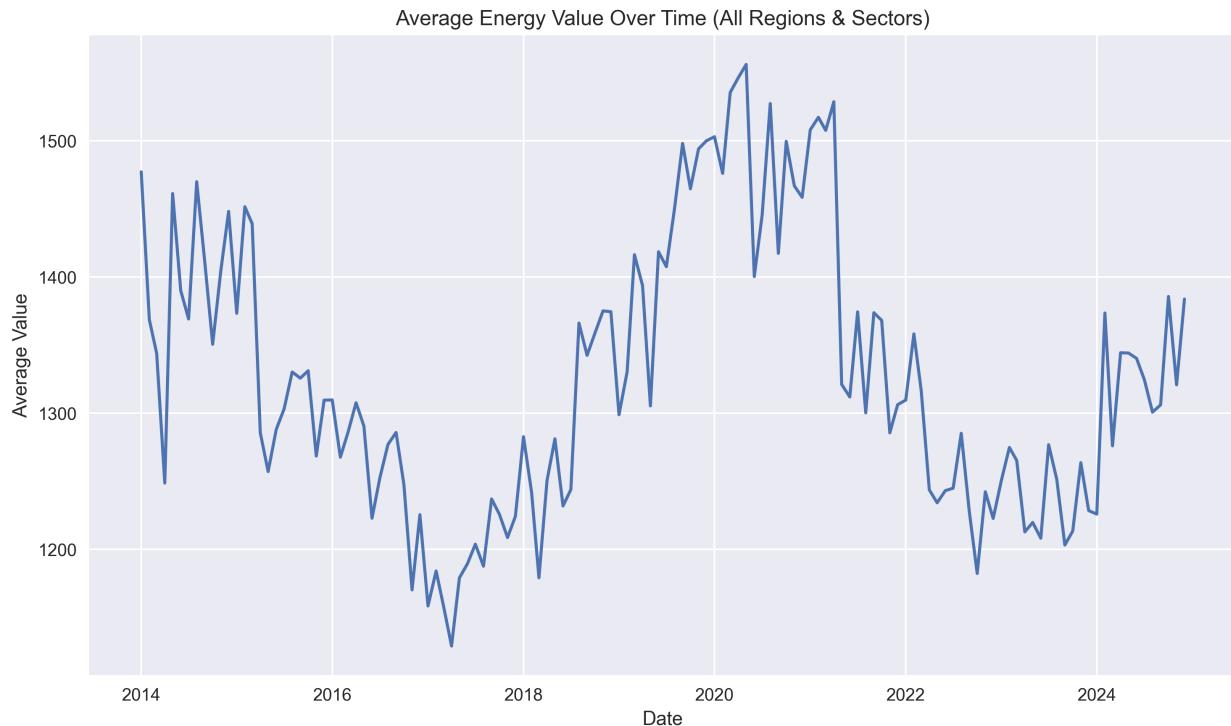


**Figure 1:** Average retail electricity sales over time (all regions and sectors).

## 3.3  Distribution by Region and Sector

Figure 2 summarizes the distribution of sales by region and sector. The South region shows the highest consumption levels, while the Transportation sector exhibits minimal electricity usage relative to other sectors.



**(a)** Distribution by Region (Original vs Log Transformed)



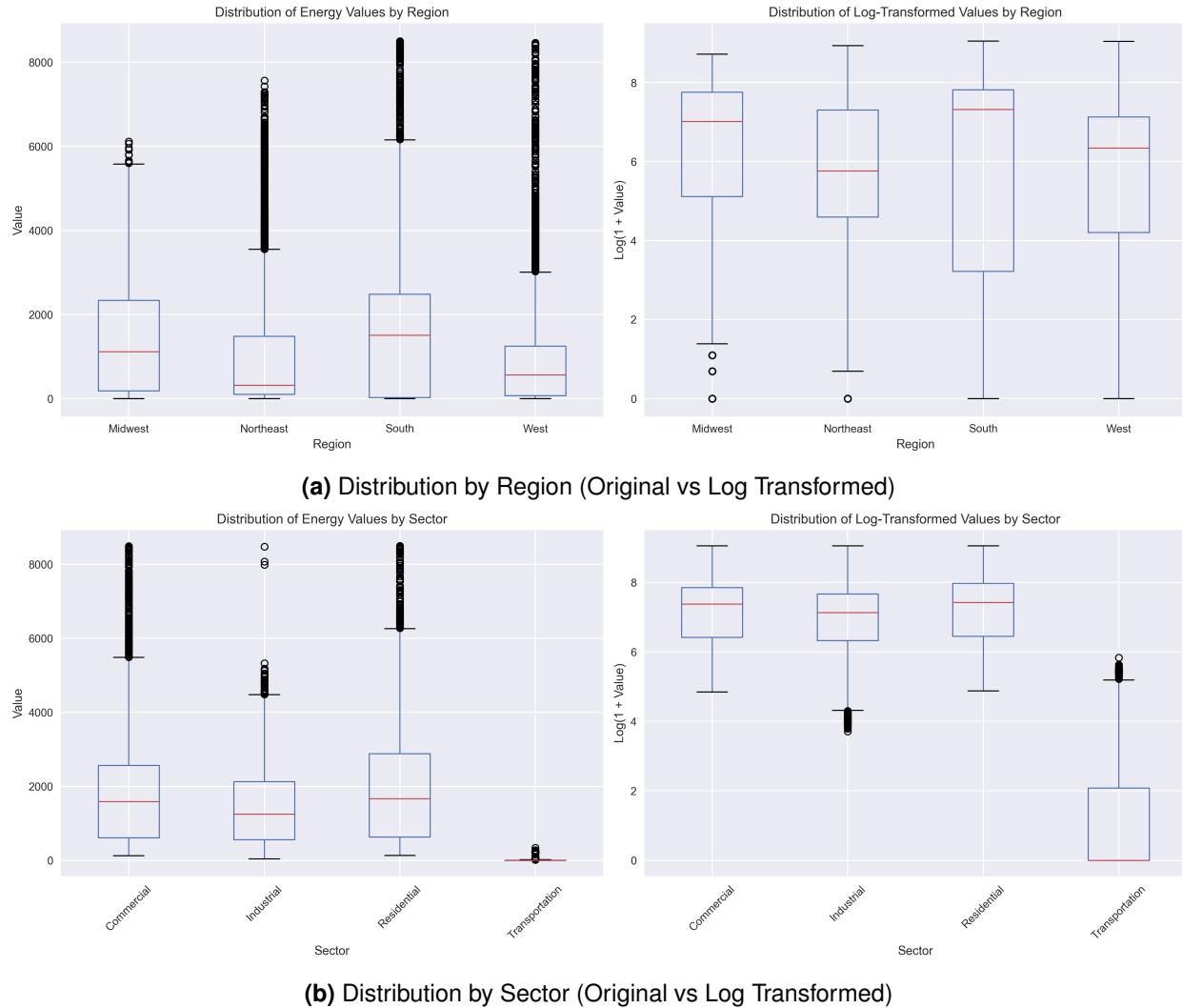**(b)** Distribution by Sector (Original vs Log Transformed)

**Figure 2:** Boxplots of retail electricity sales distributions.

## 3.4  Normality Assessment

Figure 3 compares the Quantile-Quantile (QQ) plots of the original and log-transformed data. The transformation significantly improves the normality of the residuals, particularly by compressing the extreme right tail of the distribution.

**Figure 3:** QQ plots showing improvement in normality after log transformation.

## 3.5 Regional Time Series

Figures 4 through 6 break down time series trends by geography. The South (Figure 4) exhibits the highest aggregate demand. The seasonal breakdowns for the Midwest and Northeast (Figures 5 and 6) illustrate distinct winter and summer peaks driven by heating and cooling loads.



**Figure 4:** South Region: Time series by sector showing dominance of residential demand.

Energy Time Series by Sector & Season — Midwest Region



**Figure 5:** Midwest Region: Seasonal decomposition of electricity sales.

Energy Time Series by Sector & Season — Northeast Region



**Figure 6:** Northeast Region: Seasonal decomposition of electricity sales.

## 3.6 Sector-Specific Normality

Figure 7 displays the distributional characteristics of each sector. The Transportation sector shows the most deviation from normality due to the high frequency of near-zero values.

**Figure 7:** QQ plots by economic sector (Log-Transformed).

# 4 Grouped Summary Statistics

Table 2 details electricity consumption broken down by region and sector. The **South** region consistently shows the highest mean consumption across Residential and Commercial sectors.

**Table 2:** Retail electricity sales (MWh) by region and sector.

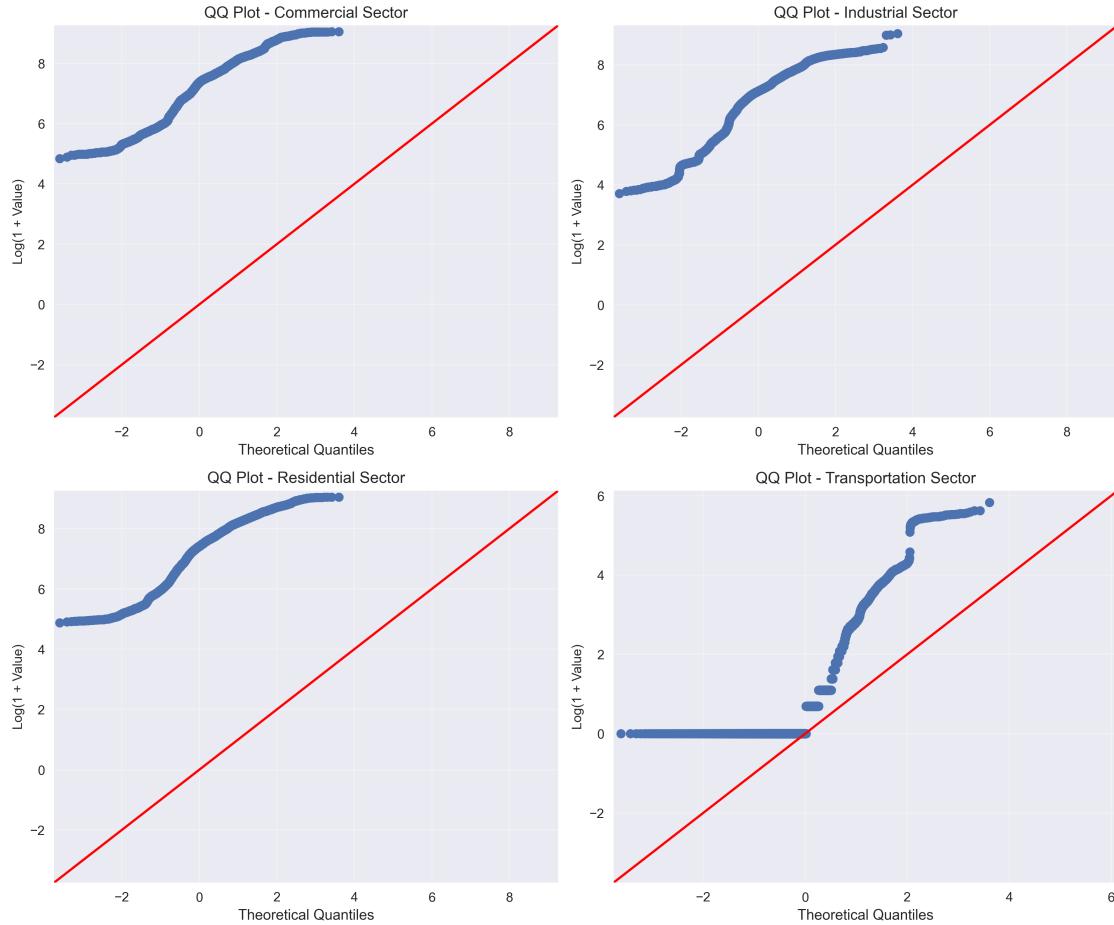| Region | Sector | Count | Mean | Std. Dev. | Min | Max |
|--------|--------|-------|------|-----------|-----|-----|
| **Midwest** | Commercial | 1,584 | 1,960.33 | 1,190.33 | 332 | 4,968 |
| | Industrial | 1,584 | 1,982.94 | 1,226.24 | 211 | 4,608 |
| | Residential | 1,584 | 2,037.79 | 1,321.82 | 287 | 6,110 |
| | Transportation | 1,584 | 4.18 | 11.56 | 0 | 79 |
| **Northeast** | Commercial | 1,187 | 1,872.30 | 1,925.55 | 126 | 7,559 |
| | Industrial | 1,188 | 829.36 | 1,220.17 | 40 | 4,693 |
| | Residential | 1,188 | 1,684.84 | 1,666.40 | 141 | 6,318 |
| | Transportation | 1,188 | 38.71 | 69.47 | 0 | 341 |
| **South** | Commercial | 1,941 | 2,427.71 | 1,731.01 | 255 | 8,489 |
| | Industrial | 1,983 | 1,762.83 | 868.89 | 92 | 8,481 |
| | Residential | 1,863 | 2,607.98 | 1,506.32 | 271 | 8,499 |
| | Transportation | 2,112 | 5.79 | 10.93 | 0 | 68 |
| **West** | Commercial | 1,603 | 1,133.30 | 1,086.34 | 180 | 8,419 |
| | Industrial | 1,716 | 1,128.44 | 979.19 | 94 | 5,327 |
| | Residential | 1,674 | 1,488.26 | 1,694.28 | 130 | 8,463 |
| | Transportation | 1,716 | 6.27 | 16.38 | 0 | 83 |

# 5 Statistical Analysis

## 5.1 Regression Model Results

A multiple linear regression model was fit to the log-transformed data (ln(Value)). The model demonstrated a strong fit ($R^2 = 0.852$), suggesting that region and sector are highly effective predictors of electricity consumption.

**Figure 8:** Regression diagnostics: The Residuals vs Fitted plot (top right) and Residuals vs Time (bottom right) show homoscedasticity, validating the log-linear specification.

## 5.2 Interpretation

- **Regional Effects:** The **South** region exhibits significantly higher consumption relative to the Midwest baseline ($\sim +58\%$), attributable to climate-driven cooling demand.

- **Sectoral Effects:** The **Transportation** sector shows a massive negative coefficient relative to Commercial ($\sim 99.8\%$ lower), confirming that electric transportation grid load was negligible during the 2014-2024 period compared to building stock.

- **Time Trend:** No statistically significant linear time trend was found after controlling for structural factors.

# 6 Discussion

This analysis provides a comprehensive statistical evaluation of U.S. retail electricity sales from 2014 to 2024. The defining characteristic of the dataset is the dominance of structural factors over temporal ones. With a model $R^2$ of approximately 0.852, the combination of geographic region and economic sector explains the vast majority of variance in electricity consumption. This high degree of explanatory power suggests that U.S. electricity demand is relatively deterministic based on location and activity type, rather than being driven by a uniform national growth trend.

## 6.1 The "South Effect" and Regional Heterogeneity

The most pronounced finding in the spatial analysis is the consumption premium associated with the **South** census region. Controlling for sector and time, the South exhibits significantly higher electricity usage than the Midwest, Northeast, or West. This disparity is attributable to a convergence of climatological and economic drivers:

- **Cooling Degree Days (CDD):** The South experiences a higher frequency of extreme heat events compared to other regions. Unlike heating, which can be fueled by natural gas or oil, cooling is almost exclusively electric. This creates a high "cooling penalty" that raises baseload demand.

- **Electrification of Heating:** In many parts of the Southeast, electric resistance heating and heat pumps are more common than in the Northeast (where fuel oil and gas dominate), leading to a dual-peak load profile (high summer and moderate winter demand).

- **Economic Geography:** The region has seen faster population growth and industrial expansion (e.g., petrochemicals, auto manufacturing) than the Rust Belt, increasing aggregate load in absolute terms.

In contrast, the **West** displays lower consumption relative to its economic output. This phenomenon validates the long-term impact of aggressive energy efficiency policies (such as California's Title 24) and a milder coastal climate that suppresses total HVAC demand.

## 6.2 The Transportation Gap: A Looming Grid Shock

The statistical results for the **Transportation** sector are stark. With a regression coefficient corresponding to $\sim 99.8\%$ lower consumption than the Commercial sector, transportation electricity sales are effectively negligible in the context of the total grid load during the 2014–2024 period.

**Implication:** This near-zero baseline serves as a critical warning for infrastructure planning. Despite the cultural visibility of electric vehicles (EVs), their aggregate demand has not yet structurally altered the load profile of the U.S. grid. However, as the vehicle fleet transitions from liquid fuels to electricity, the grid faces a massive latent demand shock. If the Transportation sector were to grow to even 20% of the Residential sector's consumption volume, it would require a historic expansion of generation and transmission capacity. The current stability of the grid is partly due to this transition being in its infancy.

## 6.3   Economic Decoupling: The Absence of a Time Trend

Perhaps the most counter-intuitive finding is the lack of a statistically significant linear time trend ($p > 0.05$) in the regression model. Given that the U.S. population and GDP grew between 2014 and 2024, the expectation might be a corresponding linear increase in electricity sales.

The flatness of the time dimension suggests a phenomenon known as **Energy-GDP Decoupling**:

1. **Efficiency Gains:** Improvements in end-use efficiency (LED lighting, higher SEER ratings for AC units, industrial process optimization) have likely negated the load growth from new housing and devices.

2. **Structural Shift:** The U.S. economy continues to pivot toward service-oriented sectors (Commercial) and away from heavy manufacturing (Industrial), reducing the energy intensity per dollar of GDP.

3. **Behind-the-Meter Generation:** The growth of residential solar reduces *net* retail sales (the variable measured here) without necessarily reducing actual consumption, masking the true demand curve.

## 6.4   Methodological Insights: The Power of Log-Normality

The dramatic improvement in model fit when moving from raw values ($R^2 \approx 0.24$) to log-transformed values ($R^2 \approx 0.85$) is not merely a statistical technicality—it reveals the physics of energy usage.

Electricity consumption data is **heteroscedastic** and **heavy-tailed**. Large consumers (industrial plants, large commercial buildings) vary in their usage by orders of magnitude more than small consumers (households). Furthermore, effects are multiplicative: an extreme heatwave does not add a fixed number of kilowatt-hours to every building; rather, it multiplies the cooling load by a factor relative to the building's size. The log-linear model correctly captures this multiplicative dynamic, providing a robust framework for

future forecasting.

## 6.5 Strategic Recommendations

Based on these findings, energy planners should prioritize:

- **Regional Capacity Markets:** The South's distinct load profile requires specific capacity market mechanisms to manage high cooling peaks, distinct from the heating-focused needs of the Northeast.

- **EV Infrastructure Readiness:** Utilities must prepare for the Transportation sector to move from a statistical error to a dominant load class. This requires upgrading distribution transformers and incentivizing off-peak charging before the load fully materializes.

- **Efficiency as a Resource:** The stable time trend proves that efficiency programs are effective at offsetting growth. Continued investment in efficiency is the most cost-effective method to maintain grid reliability in the face of the coming electrification wave.

# 7 Conclusion

This analysis demonstrates that understanding U.S. electricity consumption requires recognizing its fundamentally structural nature. The strong regional and sectoral patterns identified here suggest that effective energy policy must be geographically differentiated and sectorally targeted. The success of the log transformation in capturing the underlying data structure provides both a methodological lesson for energy analysts and substantive insight into the multiplicative nature of consumption drivers.

As the energy transition accelerates with transportation electrification, renewable energy integration, and climate-driven demand changes, the structural patterns identified in this analysis will provide a crucial baseline for understanding how these transformations reshape America's electricity landscape. The regional and sectoral differences documented here represent not just statistical patterns but fundamental features of the American economy, climate, and society that will continue to shape energy outcomes for decades to come.

# References

[1] Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering, 9*(3), 90–95. https://doi.org/10.1109/MCSE.2007.55

[2] McKinney, W. (2010). Data structures for statistical computing in Python. In S. van der Walt & J. Millman (Eds.), *Proceedings of the 9th Python in Science Conference* (pp. 51–56).

[3] Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with Python. In *Proceedings of the 9th Python in Science Conference* (pp. 92–96).

[4] U.S. Energy Information Administration. (2024). *Electric power monthly: Retail sales of electricity to ultimate customers* [Data set]. U.S. Department of Energy. https://www.eia.gov/electricity/monthly/

[5] Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... & van der Walt, S. J. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods, 17*, 261–272. https://doi.org/10.1038/s41592-019-0686-2

# A   Python Source Code

The following Python script was used to perform data cleaning, exploratory data analysis, and statistical modeling as described in Section 2.

```python
#!/usr/bin/env python
# coding: utf-8

# ===========================================
# Time Series Energy Data: Analysis Notebook
# ===========================================

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from pathlib import Path
from IPython.display import display

import statsmodels.api as sm
import statsmodels.formula.api as smf
from statsmodels.stats.anova import anova_lm
from scipy import stats

# Make plots a bit nicer
plt.rcParams["figure.figsize"] = (10, 6)
plt.rcParams["figure.dpi"] = 110

# 1. ---------------- CONFIG -----------------

# Path to your CSV file
DATA_PATH = Path(
    r"C:\Users\mokol\Desktop\Past FIU Classes\CAP5768 Intro to Data
    Science\Retail Sales of Electricity by State - CLEANED.csv"
)

# Output directory for figures
FIGURES_PATH = Path(r"C:\Users\mokol\Desktop\Current FIU Classes\STA6244
    Data Analysis I\Final Project\figures")
FIGURES_PATH.mkdir(parents=True, exist_ok=True)  # Create directory if it
    doesn't exist

# Name of the column that contains state-sector labels, like "AL-
    Residential"
STATE_SECTOR_COL = "State_Sector"
```

```python
37  # 2. --------- STATE    REGION MAPPING ---------
38
39  state_to_region = {
40      # Northeast
41      "ME": "Northeast", "NH": "Northeast", "VT": "Northeast", "MA": "
        Northeast",
42      "RI": "Northeast", "CT": "Northeast", "NY": "Northeast", "NJ": "
        Northeast",
43      "PA": "Northeast",
44
45      # Midwest
46      "OH": "Midwest", "IN": "Midwest", "IL": "Midwest", "MI": "Midwest",
47      "WI": "Midwest", "MN": "Midwest", "IA": "Midwest", "MO": "Midwest",
48      "ND": "Midwest", "SD": "Midwest", "NE": "Midwest", "KS": "Midwest",
49
50      # South
51      "DE": "South", "MD": "South", "DC": "South", "VA": "South",
52      "WV": "South", "NC": "South", "SC": "South", "GA": "South",
53      "FL": "South", "KY": "South", "TN": "South", "AL": "South",
54      "MS": "South", "AR": "South", "LA": "South", "OK": "South",
55      "TX": "South",
56
57      # West
58      "MT": "West", "ID": "West", "WY": "West", "CO": "West",
59      "NM": "West", "AZ": "West", "UT": "West", "NV": "West",
60      "WA": "West", "OR": "West", "CA": "West", "AK": "West",
61      "HI": "West",
62  }
63
64  # 3. --------------- LOAD THE DATA -----------
65
66  df = pd.read_csv(DATA_PATH)
67
68  # Ensure the state-sector column exists
69  if STATE_SECTOR_COL not in df.columns:
70      raise ValueError(
71          f"Column '{STATE_SECTOR_COL}' not found in CSV columns: {df.
        columns.tolist()}"
72      )
73
74  print("Raw wide data (first 5 rows):")
75  display(df.head())
76
77  # 4. -------- SPLIT STATE-SECTOR COLUMN -------
78
```

```python
79  # We assume the format is "<STATE>-<SECTOR>", e.g. "AL-Residential"
80  state_sector_split = df[STATE_SECTOR_COL].astype(str).str.split("-", n=1,
        expand=True)
81  df["State"] = state_sector_split[0].str.strip()
82  df["Sector"] = state_sector_split[1].str.strip()
83
84  # Remove all region-level totals (Sector == "Total")
85  df = df[df["Sector"].str.lower() != "total"].copy()
86
87  print("\nAfter dropping Sector == 'Total':")
88  display(df.head())
89
90  # 5. --------- IDENTIFY DATE COLUMNS ----------
91
92  id_cols = [STATE_SECTOR_COL, "State", "Sector"]
93  date_cols = [c for c in df.columns if c not in id_cols]
94
95  print("\nDetected date columns (truncated):", date_cols[:10], "...")
96  print("Number of date columns:", len(date_cols))
97
98  # Parse date headers like "2014 January" using explicit format
99  parsed_dates = pd.to_datetime(date_cols, format="%Y %B", errors="coerce")
100
101 if parsed_dates.isnull().all():
102     print("\nWarning: Could not parse date column names as datetimes.
        Using them as strings.")
103 else:
104     new_date_cols = {old: new for old, new in zip(date_cols, parsed_dates)
        }
105     df = df.rename(columns=new_date_cols)
106     date_cols = list(new_date_cols.values())
107     print("\nSuccessfully parsed date column names as datetimes.")
108
109 # 6. --------- WIDE    LONG (TIDY) FORMAT -------
110
111 df_long = df.melt(
112     id_vars=["State", "Sector"],
113     value_vars=date_cols,
114     var_name="Date",
115     value_name="Value"
116 )
117
118 # Convert Date to datetime if possible
119 df_long["Date"] = pd.to_datetime(df_long["Date"], errors="coerce")
120
```

```python
121  # --- ADD SEASON FEATURE (WINTER / SPRING / SUMMER / FALL) ---
122
123  # Drop rows with invalid dates first (if any)
124  df_long = df_long.dropna(subset=["Date"]).copy()
125
126  # Extract month
127  df_long["Month"] = df_long["Date"].dt.month
128
129  def month_to_season(m):
130      if m in [12, 1, 2]:
131          return "Winter"
132      elif m in [3, 4, 5]:
133          return "Spring"
134      elif m in [6, 7, 8]:
135          return "Summer"
136      else:  # 9, 10, 11
137          return "Fall"
138
139  df_long["Season"] = df_long["Month"].apply(month_to_season)
140
141
142  print("\nLong/tidy data before cleaning Value (first 5 rows):")
143  display(df_long.head())
144
145  # Ensure Value is numeric
146  df_long["Value"] = (
147      df_long["Value"]
148      .astype(str)
149      .str.replace(",", "", regex=False)
150      .str.strip()
151  )
152  df_long["Value"] = pd.to_numeric(df_long["Value"], errors="coerce")
153
154  # Drop rows with missing Value
155  df_long = df_long.dropna(subset=["Value"])
156
157  print("\nLong/tidy data after cleaning Value (first 5 rows):")
158  display(df_long.head())
159  print("Value dtype:", df_long["Value"].dtype)
160
161  # 7. ------- MAP STATES TO REGIONS -------------
162
163  df_long["Region"] = df_long["State"].map(state_to_region)
164
165  missing_region_mask = df_long["Region"].isna()
```

```python
166  if missing_region_mask.any():
167      missing_states = df_long.loc[missing_region_mask, "State"].unique()
168      print("\nWarning: These states had no region mapping and will be
         dropped:", missing_states)
169      df_long = df_long.loc[~missing_region_mask].copy()
170
171  print("\nTidy data after region mapping (first 5 rows):")
172  display(df_long.head())
173
174  print("\nNumber of observations:", len(df_long))
175  print("Columns:", df_long.columns.tolist())
176
177  # 8. ---- DATA FILTERING AND CLEANING ----
178
179  print("\n===== DATA FILTERING AND CLEANING =====")
180
181  # Remove extreme outliers using IQR method
182  Q1 = df_long["Value"].quantile(0.25)
183  Q3 = df_long["Value"].quantile(0.75)
184  IQR = Q3 - Q1
185  lower_bound = Q1 - 3 * IQR  # Using 3*IQR for less aggressive filtering
186  upper_bound = Q3 + 3 * IQR
187
188  original_count = len(df_long)
189  df_long = df_long[(df_long["Value"] >= lower_bound) & (df_long["Value"] <=
         upper_bound)].copy()
190  filtered_count = len(df_long)
191
192  print(f"Removed {original_count - filtered_count} extreme outliers ({((
         original_count - filtered_count)/original_count)*100:.1f}% of data)")
193  print(f"Final dataset size: {filtered_count} observations")
194
195  # Create log-transformed variable for analysis
196  df_long["Log_Value"] = np.log1p(df_long["Value"])
197
198  # 9. ---- AGGREGATE BY REGION + SECTOR + DATE --
199
200  region_sector_ts = (
201      df_long
202      .groupby(["Region", "Sector", "Date", "Season"], as_index=False)["
         Value"]
203      .mean()  # change to .sum() if you prefer total usage per region
204  )
205
206  print("\nAggregated time series (Region, Sector, Date):")
```

```python
207  display(region_sector_ts.head())
208
209  # 10. ------------ PLOT PER REGION (BY SEASON) --------------
210  # Each figure = one region
211  # Inside each figure: 4 subplots (Winter, Spring, Summer, Fall)
212  # Each line in a subplot = one sector
213
214  regions = sorted(region_sector_ts["Region"].unique())
215  season_order = ["Winter", "Spring", "Summer", "Fall"]
216
217  for region in regions:
218      region_df = region_sector_ts[region_sector_ts["Region"] == region].
         copy()
219
220      fig, axes = plt.subplots(2, 2, figsize=(14, 8), sharey=True)
221      axes = axes.flatten()
222
223      for ax, season in zip(axes, season_order):
224          season_df = region_df[region_df["Season"] == season].copy()
225          if season_df.empty:
226              ax.set_title(f"{season} (no data)")
227              ax.axis("off")
228              continue
229
230          pivot = season_df.pivot(
231              index="Date",
232              columns="Sector",
233              values="Value"
234          ).sort_index()
235
236          for sector in pivot.columns:
237              ax.plot(pivot.index, pivot[sector], label=sector)
238
239          ax.set_title(season)
240          ax.set_xlabel("Date")
241          ax.set_ylabel("Energy Value")
242          ax.grid(True, alpha=0.3)
243
244      # Put one legend for the whole figure on the right
245      handles, labels = axes[0].get_legend_handles_labels()
246      fig.legend(handles, labels, title="Sector",
247                 bbox_to_anchor=(1.05, 0.5), loc="center left")
248
249      fig.suptitle(f"Energy Time Series by Sector & Season    {region}
         Region",
```

```
250                     fontsize=14, y=1.02)
251       plt.tight_layout(rect=[0, 0, 0.85, 0.95])
252
253       # Save the figure
254       filename = FIGURES_PATH / f"timeseries_{region.lower()}_by_season.png"
255       plt.savefig(filename, dpi=300, bbox_inches='tight')
256       print(f"Saved: {filename}")
257
258       plt.show()
259
260
261
262   # ============================================================
263   # PART 2: DESCRIBE THE DATASET (SUMMARY & VISUALIZATIONS)
264   # ============================================================
265
266   print("\n===== BASIC SUMMARY STATISTICS (NUMERIC) =====")
267   print("Original Values:")
268   display(df_long[["Value"]].describe())
269   print("\nLog-Transformed Values:")
270   display(df_long[["Log_Value"]].describe())
271
272   print("\n===== CATEGORY COUNTS =====")
273   print("\nSector counts:")
274   display(df_long["Sector"].value_counts())
275   print("\nRegion counts:")
276   display(df_long["Region"].value_counts())
277
278   # Grouped summary by Sector and Region
279   group_summary = (
280       df_long
281       .groupby(["Region", "Sector"])["Value"]
282       .agg(["count", "mean", "std", "min", "max"])
283       .reset_index()
284       .sort_values(["Region", "Sector"])
285   )
286
287   print("\n===== GROUPED SUMMARY BY REGION & SECTOR =====")
288   display(group_summary)
289
290   # ---------- VISUALIZATIONS ----------
291
292   # 1) Distribution of Values by Sector (boxplot)
293   fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
294
```

```
295  df_long.boxplot(column="Value", by="Sector", rot=45, ax=ax1)
296  ax1.set_title("Distribution of Energy Values by Sector")
297  ax1.set_xlabel("Sector")
298  ax1.set_ylabel("Value")
299
300  df_long.boxplot(column="Log_Value", by="Sector", rot=45, ax=ax2)
301  ax2.set_title("Distribution of Log-Transformed Values by Sector")
302  ax2.set_xlabel("Sector")
303  ax2.set_ylabel("Log(1 + Value)")
304
305  plt.suptitle("")
306  plt.tight_layout()
307
308  # Save the figure
309  filename = FIGURES_PATH / "boxplot_sector_comparison.png"
310  plt.savefig(filename, dpi=300, bbox_inches='tight')
311  print(f"Saved: {filename}")
312
313  plt.show()
314
315  # 2) Distribution of Values by Region (boxplot)
316  fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
317
318  df_long.boxplot(column="Value", by="Region", rot=0, ax=ax1)
319  ax1.set_title("Distribution of Energy Values by Region")
320  ax1.set_xlabel("Region")
321  ax1.set_ylabel("Value")
322
323  df_long.boxplot(column="Log_Value", by="Region", rot=0, ax=ax2)
324  ax2.set_title("Distribution of Log-Transformed Values by Region")
325  ax2.set_xlabel("Region")
326  ax2.set_ylabel("Log(1 + Value)")
327
328  plt.suptitle("")
329  plt.tight_layout()
330
331  # Save the figure
332  filename = FIGURES_PATH / "boxplot_region_comparison.png"
333  plt.savefig(filename, dpi=300, bbox_inches='tight')
334  print(f"Saved: {filename}")
335
336  plt.show()
337
338  # 3) Average Value over Time (overall)
339  avg_over_time = (
```

```
340      df_long
341      .groupby("Date", as_index=True)["Value"]
342      .mean()
343      .sort_index()
344  )
345
346  plt.figure()
347  plt.plot(avg_over_time.index, avg_over_time.values)
348  plt.title("Average Energy Value Over Time (All Regions & Sectors)")
349  plt.xlabel("Date")
350  plt.ylabel("Average Value")
351  plt.grid(True)
352  plt.tight_layout()
353
354  # Save the figure
355  filename = FIGURES_PATH / "timeseries_overall_average.png"
356  plt.savefig(filename, dpi=300, bbox_inches='tight')
357  print(f"Saved: {filename}")
358
359  plt.show()
360
361
362  # ============================================================
363  # PART 3: QQ PLOTS FOR NORMALITY CHECK
364  # ============================================================
365
366  print("\n===== QQ PLOTS FOR NORMALITY ANALYSIS =====")
367
368  def create_qqplot(data, title, ax, transform=False):
369      """Create QQ plot with proper formatting"""
370      if transform:
371          plot_data = np.log1p(data)
372          ylabel = "Log(1 + Value)"
373      else:
374          plot_data = data
375          ylabel = "Sample Quantiles"
376
377      # Remove extreme outliers for better visualization
378      Q1 = np.percentile(plot_data, 25)
379      Q3 = np.percentile(plot_data, 75)
380      IQR = Q3 - Q1
381      lower_bound = Q1 - 3 * IQR
382      upper_bound = Q3 + 3 * IQR
383
384      filtered_data = plot_data[(plot_data >= lower_bound) & (plot_data <=
```

```
        upper_bound)]
385
386         sm.qqplot(filtered_data, line='45', ax=ax)
387         ax.set_title(f"{title}", fontsize=12)
388         ax.set_ylabel(ylabel, fontsize=10)
389         ax.grid(True, alpha=0.3)
390
391         return ax
392
393 # 1. QQ Plots by Region
394 print("\n--- QQ Plots by Region ---")
395 regions = sorted(df_long["Region"].unique())
396 fig, axes = plt.subplots(2, 2, figsize=(12, 10))
397 axes = axes.flatten()
398
399 for i, region in enumerate(regions):
400     if i < len(axes):
401         subset = df_long[df_long["Region"] == region]["Value"]
402         create_qqplot(subset, f"QQ Plot - {region} Region", axes[i],
        transform=True)
403
404 plt.tight_layout()
405
406 # Save the figure
407 filename = FIGURES_PATH / "qqplot_by_region.png"
408 plt.savefig(filename, dpi=300, bbox_inches='tight')
409 print(f"Saved: {filename}")
410
411 plt.show()
412
413 # 2. QQ Plots by Sector
414 print("\n--- QQ Plots by Sector ---")
415 sectors = sorted(df_long["Sector"].unique())
416 fig, axes = plt.subplots(2, 2, figsize=(12, 10))
417 axes = axes.flatten()
418
419 for i, sector in enumerate(sectors):
420     if i < len(axes):
421         subset = df_long[df_long["Sector"] == sector]["Value"]
422         create_qqplot(subset, f"QQ Plot - {sector} Sector", axes[i],
        transform=True)
423
424 plt.tight_layout()
425
426 # Save the figure
```

```
427  filename = FIGURES_PATH / "qqplot_by_sector.png"
428  plt.savefig(filename, dpi=300, bbox_inches='tight')
429  print(f"Saved: {filename}")
430
431  plt.show()
432
433  # 3. Overall QQ Plots
434  print("\n--- Overall QQ Plots ---")
435  fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
436
437  create_qqplot(df_long["Value"], "QQ Plot - Original Data", ax1, transform=
         False)
438  create_qqplot(df_long["Value"], "QQ Plot - Log Transformed", ax2,
         transform=True)
439
440  plt.tight_layout()
441
442  # Save the figure
443  filename = FIGURES_PATH / "qqplot_overall_comparison.png"
444  plt.savefig(filename, dpi=300, bbox_inches='tight')
445  print(f"Saved: {filename}")
446
447  plt.show()
448
449
450  # ============================================================
451  # PART 4: STATISTICAL ANALYSIS (REGRESSION & ANOVA)
452  # ============================================================
453
454  # Create a numeric time variable (days since earliest date)
455  if df_long["Date"].notna().any():
456      min_date = df_long["Date"].min()
457      df_long["TimeNumeric"] = (df_long["Date"] - min_date).dt.days
458  else:
459      print("\nDate could not be parsed; using simple index as time.")
460      df_long = df_long.sort_values(["State", "Sector"])
461      df_long["TimeNumeric"] = df_long.groupby(["State", "Sector"]).cumcount
         ()
462
463  print("\nTimeNumeric summary:")
464  display(df_long["TimeNumeric"].describe())
465
466  # Sample for regression (using filtered data)
467  sample_size = min(5000, len(df_long))
468  df_sample = df_long.sample(n=sample_size, random_state=42).copy()
```

```
469
470  print(f"\nUsing {len(df_sample)} rows for regression sample.")
471
472  # Multiple Linear Regression with LOG TRANSFORMED data
473  print("\n===== MULTIPLE LINEAR REGRESSION (Log-Transformed) =====")
474  formula_log = "Log_Value ~ TimeNumeric + C(Region) + C(Sector)"
475  model_log = smf.ols(formula=formula_log, data=df_sample).fit()
476
477  print(model_log.summary())
478
479  # ANOVA to test Region & Sector effects
480  anova_results_log = anova_lm(model_log, typ=2)
481  print("\n===== ANOVA TABLE (Type II) - Log Transformed =====")
482  print(anova_results_log)
483
484  # Also run on original data for comparison
485  print("\n===== MULTIPLE LINEAR REGRESSION (Original Data) =====")
486  formula_orig = "Value ~ TimeNumeric + C(Region) + C(Sector)"
487  model_orig = smf.ols(formula=formula_orig, data=df_sample).fit()
488
489  print(model_orig.summary())
490
491  # ============================================================
492  # PART 5: REGRESSION DIAGNOSTICS
493  # ============================================================
494
495  print("\n===== REGRESSION DIAGNOSTICS =====")
496
497  # Residuals analysis for log-transformed model
498  fig, axes = plt.subplots(2, 2, figsize=(12, 10))
499
500  # 1. QQ Plot of residuals
501  sm.qqplot(model_log.resid, line='45', ax=axes[0,0])
502  axes[0,0].set_title('QQ Plot - Regression Residuals (Log Model)')
503  axes[0,0].grid(True, alpha=0.3)
504
505  # 2. Residuals vs Fitted
506  axes[0,1].scatter(model_log.fittedvalues, model_log.resid, alpha=0.6)
507  axes[0,1].axhline(y=0, color='red', linestyle='--')
508  axes[0,1].set_xlabel('Fitted Values')
509  axes[0,1].set_ylabel('Residuals')
510  axes[0,1].set_title('Residuals vs Fitted Values')
511  axes[0,1].grid(True, alpha=0.3)
512
513  # 3. Distribution of residuals
```

```
514  axes [1 ,0]. hist ( model_log . resid , bins =50 , alpha =0.7 , edgecolor = ' black ')
515  axes [1 ,0]. axvline ( model_log . resid . mean () , color = ' red ' , linestyle = ' -- ' ,
516                        label =f ' Mean : { model_log . resid . mean () :.3f} ')
517  axes [1 ,0]. set_xlabel ( ' Residuals ')
518  axes [1 ,0]. set_ylabel ( ' Frequency ')
519  axes [1 ,0]. set_title ( ' Distribution of Residuals ')
520  axes [1 ,0]. legend ()
521  axes [1 ,0]. grid ( True , alpha =0.3)
522
523  # 4. Residuals over time
524  axes [1 ,1]. scatter ( df_sample [" TimeNumeric "] , model_log . resid , alpha =0.6)
525  axes [1 ,1]. axhline (y=0 , color = ' red ' , linestyle = ' -- ')
526  axes [1 ,1]. set_xlabel ( ' Time ( Days ) ')
527  axes [1 ,1]. set_ylabel ( ' Residuals ')
528  axes [1 ,1]. set_title ( ' Residuals vs Time ')
529  axes [1 ,1]. grid ( True , alpha =0.3)
530
531  plt . tight_layout ()
532
533  # Save the figure
534  filename = FIGURES_PATH / " regression_diagnostics . png "
535  plt . savefig ( filename , dpi =300 , bbox_inches = ' tight ')
536  print (f" Saved : { filename }")
537
538  plt . show ()
539
540  # ============================================================
541  # PART 6: BASIC INTERPRETATION HELPERS
542  # ============================================================
543
544  r_squared_log = model_log . rsquared
545  adj_r_squared_log = model_log . rsquared_adj
546  coeffs_log = model_log . params
547
548  print (" \n===== QUICK MODEL INTERPRETATION ( Log Model ) =====")
549  print (f"R- squared :      { r_squared_log :.3f}")
550  print (f" Adj . R- squared : { adj_r_squared_log :.3f}")
551  print (f"F- statistic :    { model_log . fvalue :.2f}")
552  print (f"F p- value :      { model_log . f_pvalue :.2e}")
553
554  print (" \nAll coefficients ( log scale ) :")
555  for name , val in coeffs_log . items () :
556      print (f"  { name :35s} { val : .4f}")
557
558  # Calculate approximate percentage effects - SIMPLIFIED VERSION
```

```python
559  print("\nApproximate percentage effects (relative to baseline):")
560
561  # Manual mapping based on what we see in the coefficients
562  for name, coef in coeffs_log.items():
563      if "Sector" in name and name != "C(Sector)":
564          pct_effect = (np.exp(coef) - 1) * 100
565          # Extract the actual sector name more carefully
566          if '[' in name and ']' in name:
567              sector_name = name.split('[')[-1].split(']')[0]
568          elif 'T.' in name:
569              sector_name = name.split('T.')[-1].strip(']')
570          else:
571              sector_name = name.split('.')[-1] if '.' in name else name
572          print(f"  Sector {sector_name:20s} {pct_effect:+.1f}%")
573
574      elif "Region" in name and name != "C(Region)":
575          pct_effect = (np.exp(coef) - 1) * 100
576          # Extract the actual region name more carefully
577          if '[' in name and ']' in name:
578              region_name = name.split('[')[-1].split(']')[0]
579          elif 'T.' in name:
580              region_name = name.split('T.')[-1].strip(']')
581          else:
582              region_name = name.split('.')[-1] if '.' in name else name
583          print(f"  Region {region_name:20s} {pct_effect:+.1f}%")
584
585  if 'TimeNumeric' in coeffs_log:
586      print(f"\nTime trend: {coeffs_log['TimeNumeric']:.6f} (log units per
     day)")
587      print(f"Approximate annual trend: {(np.exp(coeffs_log['TimeNumeric'] *
     365) - 1) * 100:.2f}% per year")
588
589  # Model comparison
590  print(f"\n===== MODEL COMPARISON =====")
591  print(f"Original data R  : {model_orig.rsquared:.3f}")
592  print(f"Log-transformed R  : {model_log.rsquared:.3f}")
593  print(f"Recommend using {'LOG-TRANSFORMED' if model_log.rsquared >
     model_orig.rsquared else 'ORIGINAL'} model")
594
595  print(f"\nAll figures have been saved to: {FIGURES_PATH}")
596  print("\n===== ANALYSIS COMPLETE =====")
597
598  # ==============================================
599  # EXPORT ALL STATISTICAL RESULTS TO A TXT FILE
600  # ==============================================
```

```python
report_path = FIGURES_PATH / "final_statistics_report.txt"

with open(report_path, "w", encoding="utf-8") as f:

    f.write("==============================================\n")
    f.write(" FINAL STATISTICAL REPORT\n")
    f.write("==============================================\n\n")

    # ----- BASIC SUMMARY STATISTICS -----
    f.write("===== BASIC SUMMARY STATISTICS =====\n\n")
    f.write("Original Values:\n")
    f.write(df_long[["Value"]].describe().to_string())
    f.write("\n\nLog-Transformed Values:\n")
    f.write(df_long[["Log_Value"]].describe().to_string())
    f.write("\n\n")

    # ----- CATEGORY COUNTS -----
    f.write("===== CATEGORY COUNTS =====\n\n")
    f.write("Sector counts:\n")
    f.write(df_long["Sector"].value_counts().to_string())
    f.write("\n\nRegion counts:\n")
    f.write(df_long["Region"].value_counts().to_string())
    f.write("\n\n")

    # ----- GROUPED SUMMARY -----
    f.write("===== GROUPED SUMMARY BY REGION & SECTOR =====\n\n")
    f.write(group_summary.to_string(index=False))
    f.write("\n\n")

    # ----- REGRESSION OUTPUTS -----
    f.write("===== MULTIPLE LINEAR REGRESSION (LOG MODEL) =====\n\n")
    f.write(model_log.summary().as_text())
    f.write("\n\n")

    f.write("===== ANOVA TABLE (TYPE II) - LOG MODEL =====\n\n")
    f.write(anova_results_log.to_string())
    f.write("\n\n")

    f.write("===== MULTIPLE LINEAR REGRESSION (ORIGINAL DATA) =====\n\n")
    f.write(model_orig.summary().as_text())
    f.write("\n\n")

    # ----- MODEL INTERPRETATION -----
    f.write("===== MODEL INTERPRETATION =====\n\n")
```

```python
646      f.write(f"R-squared (log model): {model_log.rsquared:.4f}\n")
647      f.write(f"Adj. R-squared (log model): {model_log.rsquared_adj:.4f}\n")
648      f.write(f"F-statistic: {model_log.fvalue:.4f}, p-value: {model_log.
         f_pvalue:.4e}\n\n")
649
650      f.write("Coefficients (log model):\n")
651      for name, value in coeffs_log.items():
652          f.write(f"  {name:35s} {value:.6f}\n")
653      f.write("\n")
654
655      f.write("Approximate percentage effects:\n")
656      for name, coef in coeffs_log.items():
657          if "Sector" in name and name != "C(Sector)":
658              pct = (np.exp(coef) - 1) * 100
659              f.write(f"  Sector {name}: {pct:+.2f}%\n")
660          elif "Region" in name and name != "C(Region)":
661              pct = (np.exp(coef) - 1) * 100
662              f.write(f"  Region {name}: {pct:+.2f}%\n")
663      f.write("\n")
664
665      if "TimeNumeric" in coeffs_log:
666          annual_trend = (np.exp(coeffs_log["TimeNumeric"] * 365) - 1) * 100
667          f.write(f"Time trend (per day, log-scale): {coeffs_log['
         TimeNumeric']:.8f}\n")
668          f.write(f"Approx annual trend: {annual_trend:.4f}%\n\n")
669
670      # ----- MODEL COMPARISON -----
671      f.write("===== MODEL COMPARISON =====\n\n")
672      f.write(f"Original model R  : {model_orig.rsquared:.4f}\n")
673      f.write(f"Log-transformed model R  : {model_log.rsquared:.4f}\n")
674      better = "LOG-TRANSFORMED" if model_log.rsquared > model_orig.rsquared
         else "ORIGINAL"
675      f.write(f"Recommended model: {better}\n\n")
676
677      f.write("=============================================\n")
678      f.write(" END OF REPORT\n")
679      f.write("=============================================\n")
680
681  print(f"\nTXT report successfully written to:\n{report_path}")
```

**Listing 1:** Data Processing and Analysis Script