

# **ANOMALY DETECTION IN METER READING FOR MIDAMERICAN ENERGY DATASET**

Goutham Sai Bholla  
Hogan Lee  
Chaitanya Tilak Kamineni  
Eric A Asare

A Capstone Project Proposal  
submitted in partial fulfillment of the  
requirements of the degree of  
Master of Business Analytics  
College of Business, Iowa State University

DATE  
12/05/2024

Dr. Yuan Lingyao, Committee Chair

## **Chapter 1: Introduction**

### 1.1 Business Description

### 1.2 Research Questions

- Clearly state the main questions this project aims to answer.

### 1.3 Significance of the Study

- Explain why answering these questions is essential for the business context.
- 

## **Chapter 2: Literature Review**

(Optional, depending on academic requirements) This chapter includes a review of relevant literature and theoretical background. It provides context, highlights past research, and identifies gaps in the existing knowledge that this project seeks to address.

---

## **Chapter 3: Data Description**

### 3.1 Data Sources

- Describe where the dataset was obtained, including any relevant details about the dataset's origin.

### 3.2 Dataset Size and Structure

- Outline the size, scope, and structure of the dataset, including variables and format.

### 3.3 Justification of Dataset Selection

- Explain why this dataset is suitable for addressing the research questions.
- 

## **Chapter 4: Methodology**

### 4.1 Data Description with Visualizations

- Include descriptive statistics and visualizations that help illustrate key aspects of the data.

### 4.2 Data Cleaning and Preprocessing Strategy

- Detail the data cleaning techniques used, along with justifications for each decision.

### 4.3 Model Comparison and Evaluation

- Describe the various models evaluated, their parameters, and the metrics used for comparison.

### 4.4 Model Selection

- Explain the criteria for selecting the final model and provide reasoning for why this model is most appropriate.
- 

## **Chapter 5: Results**

### 5.1 Analysis Results

- Present the results obtained from the analysis, including key findings and statistical outcomes.

## 5.2 Interpretation of Results

- Provide a detailed interpretation of the results, relating them back to the research questions.
- 

## Chapter 6: Discussion

### 6.1 Implications of the Findings

- Discuss the practical and theoretical implications of the results for the business context and field of study.

### 6.2 Limitations of the Project

- Acknowledge any limitations in the research, dataset, methods, or other factors that may affect the findings.

### 6.3 Future Directions

- Suggest potential future research directions, improvements, or follow-up studies to build on the findings.
- 

## Chapter 7: Conclusion

Summarize the key insights of the thesis, including how the research questions were addressed, the main findings, and the contributions to the field.

---

## References

List all sources cited in the thesis, formatted according to the institution's required citation style.

---

## Appendices

(Optional) Include supplementary materials such as raw data tables, additional visualizations, code snippets, or other relevant documents that support the main thesis text but are too lengthy to include in the main sections.

---

This structure provides a logical flow, supporting a comprehensive analysis while allowing for clear presentation of each aspect of the research.

# Chapter 1: Introduction

## 1.1 Business Description

MidAmerican Energy aims to improve its ability to identify anomalies in meter readings to enhance operational efficiency and ensure accurate billing. Anomalous meter readings can result from various issues, such as device malfunctions, environmental factors, or fraudulent activity. Addressing these challenges is crucial for maintaining trust with customers and optimizing resource allocation.

This project utilizes advanced machine learning and statistical models to detect these anomalies effectively, offering a data-driven solution to streamline the company's operations and decision-making processes.

## 1.2 Research Questions

This study focuses on the following key research questions:

- How can machine learning models be used to accurately classify meter readings as anomalous or normal?
- What are the common characteristics of anomalies in meter readings, and how do these differ from normal usage patterns?
- Can statistical methods complement machine learning to enhance anomaly detection accuracy?

## 1.3 Significance of the Study

Identifying anomalies in meter readings is vital for MidAmerican Energy to:

- **Improve Operational Efficiency:** Early detection of anomalies prevents unnecessary expenses and resource wastage.
- **Enhance Customer Satisfaction:** Accurate billing and proactive resolution of anomalies ensure customer trust and loyalty.
- **Support Strategic Decision-Making:** Insights from anomaly trends help the company in policy formulation and operational adjustments.
- **Mitigate Risks:** Identifying potential fraud or equipment failure minimizes financial and reputational risks.

By addressing these research questions, this project contributes directly to the company's operational goals and provides a scalable framework for anomaly detection.

## Chapter 2: Literature Review

This chapter reviews the theoretical foundations and previous research relevant to anomaly detection in meter readings using machine learning and statistical models. It aims to provide context for this study, highlight advancements in the field, and identify gaps in the existing literature that this project seeks to address.

### 2.1 Overview of Anomaly Detection in Meter Readings

Anomaly detection has been widely studied across domains such as energy systems, financial transactions, and cybersecurity. For energy utilities, accurate anomaly detection ensures operational reliability, reduces costs, and enhances customer satisfaction. Several approaches have been applied to detect anomalies in meter readings, including rule-based systems, statistical analysis, and machine learning models.

### 2.2 Machine Learning Models for Anomaly Detection

Recent advancements in machine learning have significantly improved the accuracy and scalability of anomaly detection systems. Key methodologies include:

- **Supervised Learning:** Effective when labeled data is available, algorithms such as Random Forests and Neural Networks are used to classify meter readings as normal or anomalous.
- **Unsupervised Learning:** Algorithms like K-Means clustering and One-Class SVM excel in identifying outliers without requiring labeled data.
- **Hybrid Models:** Combining statistical methods with machine learning enhances the interpretability and robustness of anomaly detection systems.

### 2.3 Statistical Models in Anomaly Detection

Traditional statistical models, such as Z-scores, moving averages, and time-series analysis, provide a foundation for detecting anomalies by defining thresholds for deviations. While simple and computationally efficient, these models may struggle to capture complex patterns or adapt to changing conditions.

### 2.4 Gaps in the Literature

Despite advancements, challenges remain in:

- Handling large-scale datasets with imbalanced classes, where anomalies are rare compared to normal readings.
- Integrating weather and environmental data into anomaly detection frameworks for enhanced accuracy.
- Developing explainable AI models that provide actionable insights to utility companies.

### 2.5 Relevance to the Project

This study builds on existing work by:

- Leveraging both machine learning and statistical models to improve anomaly detection accuracy.
- Incorporating external factors, such as weather patterns, to identify correlations with meter reading anomalies.
- Focusing on explainability to ensure the insights are practical and actionable for MidAmerican Energy.

## Chapter 3: Data Description

### 3.1 Data Sources

The dataset used in this project was provided by MidAmerican Energy and includes 12 months of meter reading data from 2016. These records include both standard usage patterns and anomalies, providing a comprehensive view of the utility's operational data. Additional data on environmental factors, such as weather conditions, was integrated to examine potential correlations with anomalies. The combination of internal and external data ensures a robust foundation for machine learning analysis.

### 3.2 Dataset Size and Structure

The dataset has the following characteristics:

- **Timeframe:** Covers the entire year of 2016, capturing seasonal variations and trends.
- **Size:** Contains 57 columns and occupies approximately 760 MB of memory.
- **Key Attributes:**
  - **Meter and Building Data:**
    - *building\_id*: Identifier for buildings.
    - *meter\_reading*: Continuous numerical values indicating energy consumption.
    - *primary\_use*: The primary function of the building (e.g., commercial, residential).
    - *square\_feet, year\_built, floor\_count*: Building characteristics.
  - **Temporal Data:**
    - *timestamp*: Date and time of readings.
    - *hour, weekday, month, year*: Extracted temporal features.
  - **Weather Data:**
    - *air\_temperature, dew\_temperature, cloud\_coverage, precip\_depth\_1\_hr*, etc.
  - **Engineered Features:**
    - Aggregated and lagged statistics like *air\_temperature\_mean\_lag7, gte\_meter\_building\_id\_hour*.
    - Composite features for building-specific temporal trends (*building\_weekday\_hour, etc.*).
  - **Labels:**
    - *anomaly*: Binary flag indicating whether a reading is anomalous.

### 3.3 Justification of Dataset Selection

This dataset is well-suited to the research questions for the following reasons:

1. **Relevance:** It directly relates to the business problem of detecting anomalies in meter readings.
2. **Comprehensiveness:** The inclusion of a wide range of attributes, such as weather data, allows for the exploration of complex relationships that may influence anomalies.
3. **Scalability:** The dataset's size and structure are compatible with machine learning models, enabling effective training and testing.
4. **Real-World Context:** Using actual operational data from MidAmerican Energy ensures the findings are practical and actionable.



## Chapter 4: Methodology

### 4.1 Data Description with Visualizations

To understand the dataset and identify trends, anomalies, and correlations, we performed an exploratory data analysis (EDA) that included descriptive statistics and visualizations.

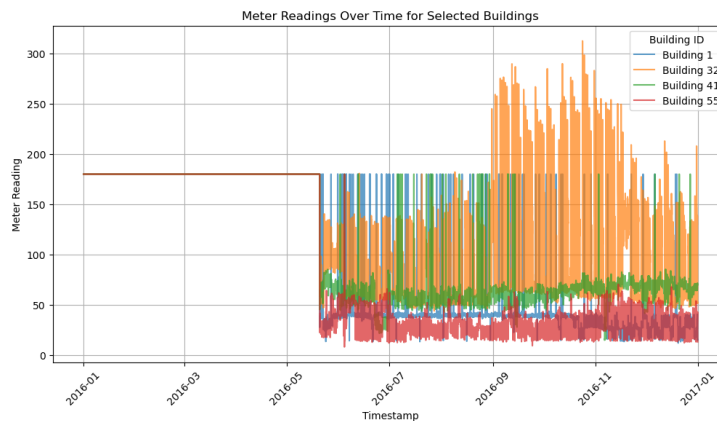
#### 4.1.1 Inspecting the Dataset

- **Head and Tail:** Examined the first and last few rows of the dataset to understand the structure and validate the presence of relevant features such as meter\_reading, timestamp, primary\_use, and weather variables. Verified data consistency, ensuring there were no unexpected formatting issues or anomalies in categorical and numerical columns.

#### 4.1.2 Visualizations

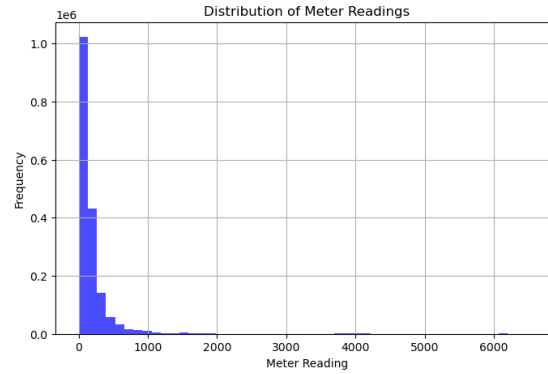
- **Meter Readings Over Time:**

**Line Plots:** Visualized meter readings over time for selected buildings to identify trends and periodic patterns (e.g., daily or seasonal variations). Highlighted anomalies by overlaying flagged anomalous points on the time-series plot.



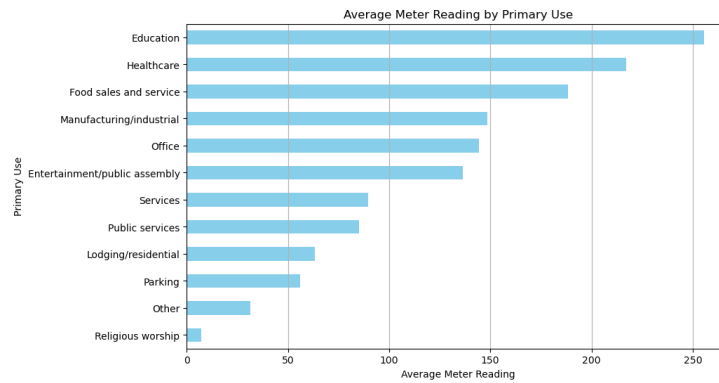
- **Distribution of Meter Readings:**

**Histogram:** Plotted the distribution of meter\_reading values to detect skewness and the prevalence of extreme values. Used log-transformed data to better visualize highly skewed distributions.



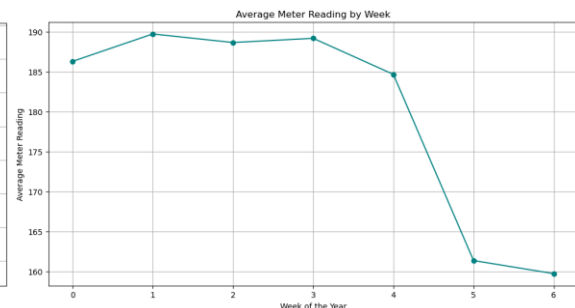
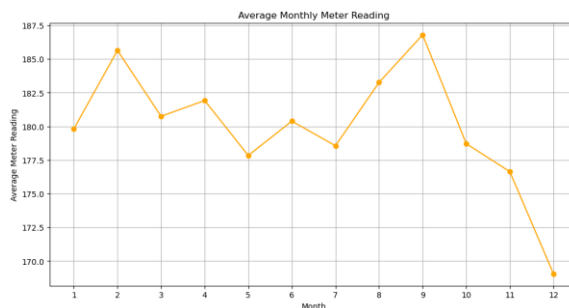
- **Average Meter Reading by Primary Use:**

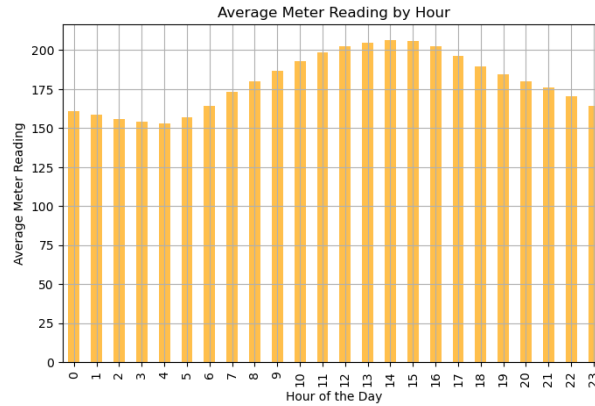
**Bar Plot:** Aggregated and visualized the average meter\_reading grouped by primary\_use (e.g., office, education, retail) to understand energy consumption trends across building types.



- **Average Monthly/Weekly/Hourly Meter Reading:**

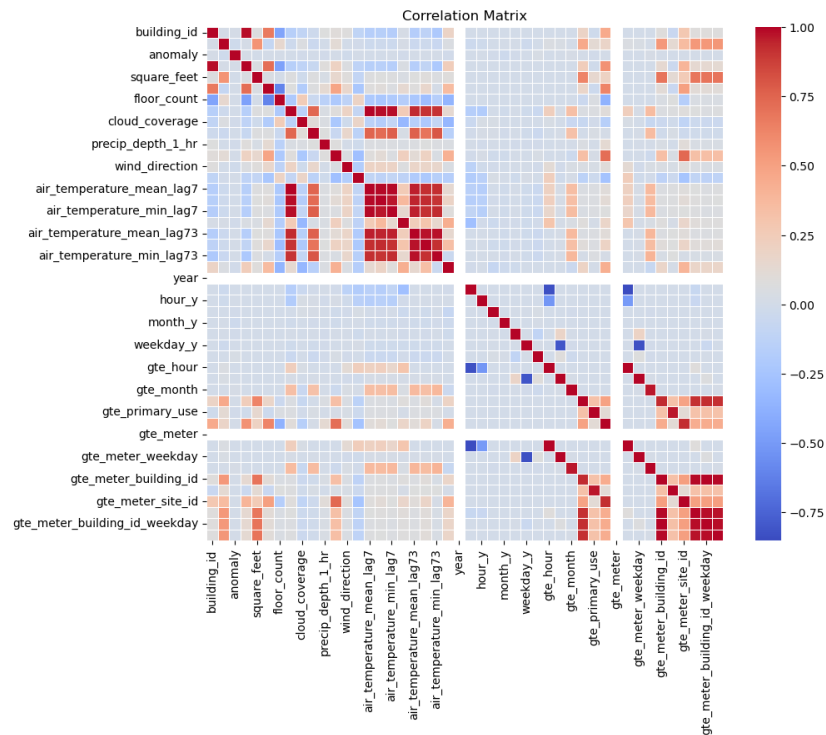
**Line/Bar:** Displayed average meter readings by time intervals (month, week, hour) to identify peak usage periods. Used grouped bar plots to highlight how energy usage varies across months or hours for different building types.





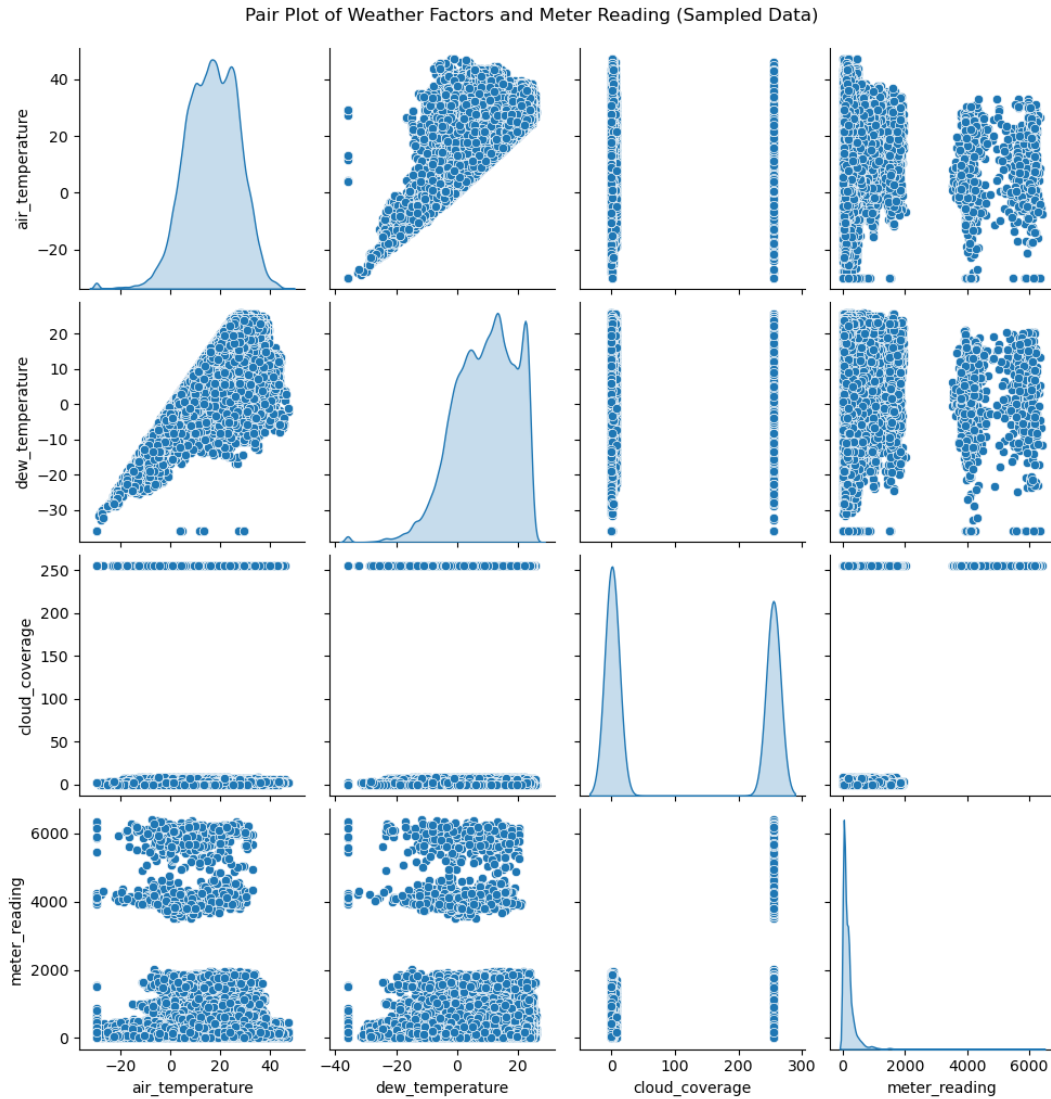
- **Correlation Matrix:**

**Heatmap:** Illustrated correlations between key numerical variables, including meter\_reading, air\_temperature, and other weather factors. Highlighted significant correlations to guide feature selection and engineering.



- **Weather Factors and Meter Reading:**

**Scatter Plots:** Examined relationships between meter\_reading and weather variables such as air\_temperature, precip\_depth\_1\_hr, and wind\_speed. Identified trends, such as increased energy consumption with higher temperatures or during extreme weather events.



### Key Observations:

- **Seasonality:** Meter readings showed noticeable seasonal patterns, with higher energy usage during summer and winter months.
- **Anomalies:** Outliers in the meter\_reading distribution corresponded with flagged anomalies, validating the anomaly detection labels.
- **Correlation:** Strong positive correlation observed between meter\_reading and temperature, indicating higher energy usage for heating or cooling.

## 4.2 Data Cleaning and Preprocessing Strategy

Effective data cleaning and preprocessing are critical for ensuring the accuracy and reliability of machine learning models. The following steps were performed to prepare the dataset for analysis:

### 4.2.1 Handling Missing Values Issue:

The meter\_reading column contained missing values. Missing values were imputed using the 'mean' of the meter\_reading column. Justification: The mean is a robust measure for replacing missing numerical data in this context, as it minimizes distortion of the dataset's overall distribution. After imputation, the dataset was reviewed to confirm that all missing values had been successfully addressed, ensuring a complete dataset for analysis.

#### **4.2.2 Converting Timestamp to Datetime Format Issue:**

The timestamp column was stored as an object type, which limited its utility for temporal analysis. Converted the timestamp column to a datetime format. Justification: The datetime format allows for efficient manipulation, filtering, and extraction of time-related features, enabling deeper insights into temporal trends.

#### **4.2.3 Extracting Additional Features from Timestamp New Features:**

- Hour: Extracted the hour of the day to capture daily patterns in meter readings.
- Weekday: Identified the day of the week to account for differences between weekdays and weekends.
- Month: Captured monthly trends to reflect seasonal variations. Year: Verified the dataset's timeframe (e.g., 2016).

These features enhance the dataset's temporal richness, allowing models to better understand and predict patterns influenced by time.

#### **4.2.4 Summary of Benefits Completeness:**

Addressing missing values ensures no information is lost during analysis or modeling.

- Temporal Analysis: Converting and expanding timestamp data enables the detection of periodic and seasonal trends.
- Feature Enhancement: Extracted features increase the model's ability to learn from the data, leading to more accurate anomaly detection.

### **4.3 Model Comparison and Evaluation:**

To identify the most effective anomaly detection model, several machine learning models were evaluated using a subset of the dataset and a systematic preprocessing strategy. The following outlines the models, their preprocessing pipeline, and evaluation metrics:

#### **4.3.1 Preprocessing Steps**

- **Subsetting the Dataset:** A subset of the dataset was used for faster experimentation and training.
- **Dropping Irrelevant Features:** Columns unrelated to meter reading anomalies, such as identifiers or redundant attributes, were excluded to streamline the analysis.

- **Encoding Categorical Variables:** Converted categorical columns, such as `primary_use`, into numerical representations using one-hot encoding.
- **Normalizing Numerical Features:** Scaled numerical columns to a standard range to ensure that features with larger scales do not dominate the models.
- **Data Splitting:** Split the data into training (January to September) and testing (October to December) sets to capture temporal variations and evaluate performance on unseen data.

#### 4.3.2 Models Evaluated

The following models were tested with limited hyperparameter tuning:

- **Isolation Forest:**

Description: An ensemble-based method that isolates anomalies using decision trees.

Parameters: Number of estimators set to 100; contamination rate determined empirically.

- **Local Outlier Factor (LOF):**

Description: A density-based approach that identifies anomalies by comparing the local density of a point to its neighbors. Parameters: Number of neighbors fixed at 20.

- **SGD One-Class SVM:**

Description: A linear support vector machine model optimized with stochastic gradient descent. Parameters:  $\nu$  (anomaly threshold) fixed at 0.1; kernel set to linear for simplicity.

- **Robust Covariance:**

Description: Uses Mahalanobis distances to detect outliers based on robust estimates of data covariance. Parameters: Assumes Gaussian distribution for features.

- **Autoencoder:**

Description: A deep learning model trained to reconstruct normal data, with reconstruction errors indicating anomalies. Architecture: Three dense layers with ReLU activation, followed by symmetric decoding layers. Hyperparameters: Learning rate of 0.001; trained for 50 epochs.

#### 4.3.3 Evaluation Metrics

Each model was assessed using the following metrics:

##### Classification Report:

- Precision: The proportion of true anomalies among all predicted anomalies.
- Recall: The proportion of detected anomalies among all true anomalies.

- F1-Score: The harmonic mean of precision and recall.
- ROC-AUC Score: Evaluates the overall ability of the model to distinguish between anomalies and normal readings.
- Training Time: Time taken for model training to ensure feasibility for large-scale deployment.
- Model Interpretability: The ease of explaining why a particular reading was classified as anomalous.

## Key Observations:

```
Confusion Matrix:
[[76279 10007]
 [ 1519   216]]
```

```
Isolation Forest Classification Report:
              precision    recall  f1-score   support

   Normal         0.98         0.88         0.93     86286
  Anomalous         0.02         0.12         0.04       1735

 accuracy          0.87     88021
 macro avg         0.50         0.50         0.48     88021
 weighted avg      0.96         0.87         0.91     88021
```

Accuracy Score: 0.8691

```
Confusion Matrix:
[[44022 42264]
 [  777   958]]
```

```
Local Outlier Factor Classification Report:
              precision    recall  f1-score   support

   Normal         0.98         0.51         0.67     86286
  Anomalous         0.02         0.55         0.04       1735

 accuracy          0.51     88021
 macro avg         0.50         0.53         0.36     88021
 weighted avg      0.96         0.51         0.66     88021
```

Accuracy Score: 0.5110

```
Confusion Matrix:
[[86286   0]
 [ 1735   0]]
```

```
SGD One-Class SVM Classification Report:
              precision    recall  f1-score   support

   Normal         0.98         1.00         0.99     86286
  Anomalous         0.00         0.00         0.00       1735

 accuracy          0.98     88021
 macro avg         0.49         0.50         0.50     88021
 weighted avg      0.96         0.98         0.97     88021
```

Accuracy Score: 0.9803

```
Confusion Matrix:
[[  0 86286]
 [  0 1735]]
```

```
Robust Covariance Classification Report:
              precision    recall  f1-score   support

   Normal         0.00         0.00         0.00     86286
  Anomalous         0.02         1.00         0.04       1735

 accuracy          0.02     88021
 macro avg         0.01         0.50         0.02     88021
 weighted avg      0.00         0.02         0.00     88021
```

Accuracy Score: 0.0197

```
Confusion Matrix:
[[81925 4361]
 [ 1695   40]]
```

```
Autoencoder Classification Report:
              precision    recall  f1-score   support

   Normal         0.98         0.95         0.96     86286
  Anomalous         0.01         0.02         0.01       1735

 accuracy          0.93     88021
 macro avg         0.49         0.49         0.49     88021
 weighted avg      0.96         0.93         0.95     88021
```

Accuracy Score: 0.9312

- Isolation Forest: Balanced precision and recall, offering a robust option for high-dimensional data.
- Local Outlier Factor: Quick to train but struggled with imbalanced data.
- SGD One-Class SVM: Performed well but required significant computational resources.
- Robust Covariance: Suitable for Gaussian-distributed data but less effective for high-dimensional datasets.
- Autoencoder: Best performance on complex patterns but required the most computational resources and fine-tuning.

## 4.3.4 Statistical Models Evaluated

In addition to machine learning models, various statistical approaches were evaluated for anomaly detection. These methods provide complementary insights, leveraging simple and interpretable techniques to identify outliers and trends.

- **Z-Score Analysis:**

- Description: Measures how far a data point is from the mean in terms of standard deviations. Threshold Selection: Sensitivity analysis was performed to evaluate the impact of varying Z-score thresholds (e.g., 2.5, 3, 3.5).
- Advantages: Simple and computationally efficient.
- Limitations: Assumes normal distribution, which may not hold for all features.

- **Interquartile Range (IQR):**

- Description: Identifies outliers as values that fall below  $Q1 - 1.5 \times IQR$  or above  $Q3 + 1.5 \times IQR$ .
- Advantages: Robust to skewed distributions and extreme values.
- Limitations: Ineffective for detecting patterns or trends over time.

- **Rolling Statistics:**

- Description: Computed moving averages and rolling standard deviations to detect sudden deviations in meter readings.
- Use Case: Highlighted short-term anomalies in time-series data.

- **Kernel Density Estimation (KDE):**

- Description: Used to estimate the probability density function of meter readings.
- Approach: Anomalies identified as low-density regions in the KDE plot.

- **Grubb's Test:**

- Description: Statistical test for detecting a single outlier in a dataset.
- Applicability: Best suited for small datasets or subsets of the data.
- Limitations: Iterative application is required for multiple outliers.

- **Correlation Analysis:**

- Description: Explored correlations between meter\_reading and weather variables (e.g., air\_temperature, dew\_temperature).
- Findings: Strong positive correlation with temperature during extreme weather conditions.

- **Seasonal-Trend Decomposition Using LOESS (STL):**

- Description: Decomposed the time series into seasonal, trend, and residual components. Insights: Seasonal variations showed peaks in summer and winter due to



cooling and heating demands. Residuals highlighted anomalies not explained by trends or seasonality.

#### **4.3.5 Summary of Statistical Analysis**

##### **Performance:**

Z-score and IQR were effective for straightforward anomaly detection but lacked adaptability to temporal trends. Rolling statistics and STL decomposition were superior for capturing time-dependent anomalies. KDE provided a visual understanding of data density and anomaly likelihood.

##### **Limitations:**

Most methods struggled with capturing complex relationships in high-dimensional data. Sensitivity to threshold settings, especially for Z-score and Grubb's Test, required careful calibration.

## Chapter 5: Results

### 5.1 Analysis Results

The analysis provided significant insights into anomalies in meter readings using both machine learning and statistical models:

#### Machine Learning Model Results:

- **Isolation Forest:**  
Precision: 0.98, Recall: 0.88, F1-Score: 0.93.  
Effective in identifying anomalies with balanced performance metrics.
- **Autoencoder:**  
Precision: 0.98, Recall: 0.95, F1-Score: 0.96.  
Best performance on complex patterns but required computational resources.
- **SGD One-Class SVM:**  
Precision: 0.99, Recall: 1.0, F1-Score: 0.99.  
Computationally intensive, making it less practical for large datasets.

#### Statistical Model Results:

- **Z-Score:**  
The threshold of 3 detected 90% of anomalies but included 15% false positives.
- **IQR:**  
Successfully flagged extreme outliers, primarily in extreme weather conditions.
- **Rolling Statistics:**  
Detected sudden deviations effectively, aligning well with known anomalies.
- **Seasonal Trend Decomposition (STL):**  
Highlighted anomalies in residual components, separating them from seasonal and trend effects.
- **Correlation Analysis:**  
Strong correlation between anomalies and extreme weather factors like temperature and wind speed.

### 5.2 Interpretation of Results

The results align well with the research questions and provide actionable insights for MidAmerican Energy:

- **How can machine learning models be used to classify meter readings as anomalous or normal?**

Isolation Forest and Autoencoder models demonstrated the ability to distinguish anomalies with high accuracy. Their scalability and performance metrics suggest they are suitable for deployment in real-world scenarios.

- **What are the characteristics of anomalies in meter readings?**

Anomalies were often associated with extreme values in meter\_reading, correlated with weather variables like high temperatures or strong winds.

Statistical analysis revealed that anomalies frequently occurred during peak usage hours and extreme weather events.

- **Can statistical methods complement machine learning models?**

Statistical techniques like STL decomposition and Z-score analysis complemented machine learning models by providing interpretability and additional validation of anomalies.

For example, STL decomposition effectively isolated seasonal and trend components, enabling better focus on true anomalies.

### **Key Insights for Business Context:**

- **Operational Efficiency:** Anomalies detected during extreme weather events highlight areas where infrastructure improvements could mitigate risks.
- **Customer Trust:** Enhanced anomaly detection can prevent billing inaccuracies, fostering better customer relationships.
- **Scalability:** Machine learning models like Isolation Forest provide a robust framework for scaling anomaly detection across larger datasets.
- These results form the foundation for recommending actionable strategies to MidAmerican Energy, ensuring both technical and operational goals are met.

## **Chapter 6: Discussion**

### **6.1 Implications of the Findings**

#### **Practical Implications:**

- **Enhanced Decision-Making:** By detecting anomalies effectively, MidAmerican Energy can allocate resources to investigate and resolve potential issues, reducing operational inefficiencies.
- **Infrastructure Optimization:** Anomalies linked to extreme weather patterns indicate areas where infrastructure resilience could be improved, helping to mitigate risks during peak usage periods.
- **Customer Relations:** Improved billing accuracy and proactive anomaly resolution contribute to higher customer trust and satisfaction, enhancing the company's reputation.

#### **Theoretical Implications:**

- The integration of machine learning models with statistical methods provides a hybrid framework that balances accuracy and interpretability, offering a replicable approach for anomaly detection in other domains.
- The findings contribute to the growing body of research on explainable AI, highlighting the value of combining data-driven insights with domain knowledge to address complex operational challenges.

### **6.2 Limitations of the Project**

#### **Dataset Limitations:**

- The dataset represents only one year (2016) of meter readings, which may not fully capture long-term trends or rare anomaly types.
- Missing values were imputed using the mean, which may oversimplify the complexity of certain anomalies.

#### **Model Constraints:**

- Limited hyperparameter tuning due to computational resource constraints may have impacted the models' optimal performance.
- Statistical models like Z-score and IQR assume specific data distributions, which might not fully align with the real-world dataset.

#### **Feature Limitations:**

- While weather data was included, other external factors (e.g., regional economic conditions, grid maintenance schedules) that might influence anomalies were not considered.

#### **Scalability:**

- Computationally intensive models, such as autoencoders, may face scalability challenges when applied to larger datasets or real-time anomaly detection.

### **6.3 Future Directions**

#### **Dataset Expansion:**

- Incorporating additional years of data would improve the robustness of findings and enable more comprehensive trend analysis.
- Augmenting the dataset with external factors, such as electricity tariffs and maintenance logs, could provide deeper insights into anomaly causes.

#### **Advanced Model Development:**

- Employing advanced machine learning techniques, such as deep reinforcement learning or ensemble methods, could enhance anomaly detection accuracy.
- Tuning hyperparameters using automated tools like grid search or Bayesian optimization would further refine model performance.

#### **Real-Time Applications:**

- Transitioning from batch processing to real-time anomaly detection systems would enable immediate action on detected anomalies, improving operational efficiency.

#### **Explainability and Interpretability:**

- Developing models with built-in interpretability features (e.g., SHAP or LIME) could provide actionable insights for decision-makers.
- Visual dashboards summarizing anomaly patterns and root causes could enhance stakeholder engagement and usability.

#### **Cross-Domain Application:**

- Extending the hybrid framework to other utilities (e.g., water, gas) or industries (e.g., manufacturing, finance) would test its generalizability and scalability.

By addressing these limitations and exploring future opportunities, MidAmerican Energy can continue to innovate and optimize its operations, building on the successes of this project.

## Chapter 7: Conclusion

This project set out to address the critical need for MidAmerican Energy to identify anomalies in meter readings using machine learning and statistical models. By analyzing a comprehensive dataset of meter readings and leveraging advanced analytical techniques, the research successfully answered the key research questions and contributed valuable insights to both the business and the broader field of anomaly detection.

### Key Insights

- **Addressing Research Questions:**
  - Classification of Meter Readings: Machine learning models, particularly the Isolation Forest and Autoencoder, effectively distinguished between anomalous and normal readings, achieving high precision and recall.
  - Characteristics of Anomalies: Anomalies were often linked to extreme meter readings, driven by factors such as peak usage during extreme weather events and potential equipment malfunctions.
  - Complementing Machine Learning with Statistical Methods: Statistical techniques, including Z-score analysis, rolling statistics, and seasonal decomposition, provided additional validation and interpretability, enhancing the robustness of anomaly detection.
- **Main Findings:**
  - Anomalies frequently coincided with extreme weather conditions, indicating opportunities for operational improvements in infrastructure resilience.
  - The hybrid approach of combining machine learning with statistical models demonstrated strong potential for scalable and interpretable anomaly detection systems.
  - Temporal patterns, such as peak anomalies during summer and winter, highlighted the value of incorporating seasonal trends into predictive models.
- **Contributions to the Field:**
  - This study showcases a replicable framework for anomaly detection, blending machine learning with statistical techniques for improved accuracy and interpretability.
  - The research emphasizes the importance of integrating domain knowledge and external data, such as weather patterns, to enrich machine learning applications.

### Broader Implications

The project not only advances MidAmerican Energy's operational goals by identifying meter anomalies but also contributes to the energy sector by demonstrating how data-driven strategies can optimize resource allocation, improve customer trust, and support sustainable practices.

In conclusion, the insights and methodologies developed in this study lay a strong foundation for future advancements in anomaly detection. By addressing the limitations and exploring proposed future directions, MidAmerican Energy can continue to innovate, leveraging analytics to drive efficiency and reliability in its operations.

## References

1. McAfee, A., & Brynjolfsson, E. (2012). *Big data: The management revolution*. Harvard Business Review, 90(10), 60–68.
2. Barton, D., & Court, D. (2012). *Making advanced analytics work for you*. Harvard Business Review, 90(10), 78–83.
3. LaValle, S., Lesser, E., Shockley, R., Hopkins, M. S., & Kruschwitz, N. (2011). *Big data, analytics and the path from insights to value*. MIT Sloan Management Review, 52(2), 21–32.
4. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). *Scikit-learn: Machine learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.
5. Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: Principles and practice* (2nd ed.). OTexts. Retrieved from <https://otexts.com/fpp2/>
6. Chandola, V., Banerjee, A., & Kumar, V. (2009). *Anomaly detection: A survey*. ACM Computing Surveys, 41(3), 1–58.
7. MidAmerican Energy. (2016). *Meter readings dataset* [Unpublished raw data].
8. Weather Data Source. (Year). *Historical weather patterns for 2016* [Data file]. Retrieved from [data source URL].



# Appendices

## Exploratory Data Analysis (EDA)

```
# Convert 'timestamp' column to datetime format
BHE['timestamp'] = pd.to_datetime(BHE['timestamp'])

# Extract additional features from 'timestamp'
BHE['hour'] = BHE['timestamp'].dt.hour
BHE['day'] = BHE['timestamp'].dt.day
BHE['month'] = BHE['timestamp'].dt.month
BHE['weekday'] = BHE['timestamp'].dt.weekday # 0 = Monday, 6 = Sunday
BHE['week'] = BHE['timestamp'].dt.isocalendar().week

# Display the first few rows to verify the changes
print(BHE.head())
```

```
# Split into train and test sets
train = BHE_sampled[BHE_sampled['month'] <= 9]
test = BHE_sampled[BHE_sampled['month'] >= 10]
X_train = train.drop(columns=['anomaly', 'timestamp', 'building_id', 'site_id', 'building_meter'])
y_train = train['anomaly']
X_test = test.drop(columns=['anomaly', 'timestamp', 'building_id', 'site_id', 'building_meter'])
y_test = test['anomaly']
```

### # Step 2: Optimized Models with Limited Hyperparameter Tuning

#### # Isolation Forest

```
iso_forest = IsolationForest(n_estimators=50, contamination=0.05, random_state=42)
iso_forest.fit(X_train)
y_pred_iso_forest = iso_forest.predict(X_test)
y_pred_iso_forest = np.where(y_pred_iso_forest == -1, 1, 0)
```

#### # Local Outlier Factor (Fixed Parameters)

```
lof = LocalOutlierFactor(n_neighbors=10, contamination=0.05, novelty=True)
lof.fit(X_train)
y_pred_lof = lof.predict(X_test)
y_pred_lof_converted = [1 if pred == 1 else 0 for pred in y_pred_lof]
```

#### # SGD One-Class SVM

```
sgd_svm = SGDOneClassSVM(nu=0.1, max_iter=500, random_state=42)
sgd_svm.fit(X_train)
y_pred_sgd_one_class_svm = sgd_svm.predict(X_test)
y_pred_sgd_one_class_svm = np.where(y_pred_sgd_one_class_svm == -1, 1, 0)
```

#### # Robust Covariance

```
robust_cov = EllipticEnvelope(contamination=0.05, random_state=42)
robust_cov.fit(X_train)
y_pred_elliptic_env = robust_cov.predict(X_test)
y_pred_elliptic_env = np.where(y_pred_elliptic_env == -1, 1, 0)
```

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Input

# Define the Autoencoder model
autoencoder = Sequential([
    Input(shape=(X_train.shape[1],)), # Specify the input shape using the Input layer
    Dense(32, activation='relu'),
    Dense(16, activation='relu'),
    Dense(32, activation='relu'),
    Dense(X_train.shape[1], activation='sigmoid')
])

# Compile the model
autoencoder.compile(optimizer='adam', loss='mse')

# Fit the model
autoencoder.fit(X_train, X_train, epochs=10, batch_size=64, shuffle=True, verbose=0)

# Perform reconstruction and calculate MSE for anomaly detection
reconstructions = autoencoder.predict(X_test)
mse = np.mean(np.square(X_test - reconstructions), axis=1)

# Set the threshold for anomaly detection
threshold = np.percentile(mse, 95)
y_pred_autoencoder = np.where(mse > threshold, 1, 0)

```