# Electricity Demand Forecasting and Analysis Report

## 1. Introduction

This report presents the analysis and forecasting of electricity demand across multiple cities in the United States. The project combines weather data, temporal features, and demand patterns to create robust forecasting models and identify patterns through clustering.

## 2. Data Collection and Preprocessing

### 2.1 Data Sources

- Weather data from 10 major cities (Dallas, Houston, LA, NYC, Philadelphia, Phoenix, San Antonio, San Diego, San Jose, Seattle)
- Electricity demand data from multiple regions
- Temporal data (date, time, season)

### 2.2 Data Preprocessing Steps

1. **Data Integration**

   - Combined weather data from multiple cities
   - Merged demand data with weather data
   - Standardized datetime formats

2. **Feature Engineering**

   - Temporal features:
     - Hour of day
     - Day of week
     - Month
     - Season (winter, spring, summer, fall)
   - Weather features:
     - Temperature
     - Humidity

- Wind speed
- Precipitation
- UV index
- Pressure

3. **Anomaly Detection**

- Z-score method (threshold: ±3)
- IQR method (1.5 * IQR)
- Isolation Forest (contamination: 0.01)

# 3. Methodology

## 3.1 Clustering Analysis

- **Method**: PCA + K-means Clustering
- **Features Used**:
  - Weather metrics
  - Temporal features
  - Demand patterns
  - Anomaly indicators
- **Implementation**:
  - Standardization of features
  - PCA for dimensionality reduction
  - K-means clustering (k=4)
- **Results**:
  - Silhouette Score: 0.413
  - Cluster Characteristics:

| Cluster | Avg Demand | Avg Temp | Avg Humidity | Avg Hour |
|---------|------------|----------|--------------|----------|
| 0 | 10088.54 | 0.707 | 0.599 | 12.75 |
| 1 | 6270.17 | 0.539 | 0.640 | 11.34 |
| 2 | 3673.19 | 0.358 | 0.662 | 10.67 |
| 3 | 5458.81 | 0.484 | 0.873 | 11.45 |

## 3.2 Forecasting Models

### 3.2.1 Original Models

1. **Linear Regression**

- Basic implementation
- Features: weather and temporal data
- Performance: MAE: 4482.21, RMSE: 4618.91

2. **XGBoost**

- Parameters:
  - n_estimators: 100
  - random_state: 42
- Performance: MAE: 3843.48, RMSE: 4141.59

3. **Random Forest**

- Parameters:
  - n_estimators: 100
  - random_state: 42
- Performance: MAE: 3690.29, RMSE: 4018.50

**3.2.2 Improved Models**

1. **Enhanced Linear Regression**

- Added polynomial features (degree=2)
- Ridge regularization (alpha=1.0)
- Standardization of features
- Performance: MAE: 40.81, RMSE: 70.30, $R^2$: 0.9997

2. **Enhanced XGBoost**

- Parameters:
  - n_estimators: 200
  - max_depth: 5
  - learning_rate: 0.1
  - subsample: 0.8
  - colsample_bytree: 0.8
- Performance: MAE: 24.77, RMSE: 35.82, $R^2$: 0.9999

3. **Enhanced Random Forest**

- Parameters:
  - n_estimators: 200
  - max_depth: 10
  - min_samples_split: 5
  - min_samples_leaf: 2
- Performance: MAE: 516.77, RMSE: 748.47, $R^2$: 0.9672

4. **LSTM Neural Network**

- Architecture:
    - Sequential model with multiple LSTM layers
    - Input shape: (look_back, n_features)
    - LSTM layers with return_sequences=True
    - Dropout layers (0.2) for regularization
    - Dense output layer
- Training:
    - Optimizer: Adam
    - Loss function: Mean Squared Error
    - Batch size: 32
    - Epochs: 50
    - Early stopping with patience=10
- Features:
    - Time series data with look-back window
    - Weather features
    - Temporal features
- Performance:
    - MAE: 15.23
    - RMSE: 22.45
    - $R^2$: 0.9999
    - Best performing model overall

# 4. Results and Discussion

## 4.1 Model Performance Comparison

| Model | MAE | RMSE | $R^2$ Score |
|---|---|---|---|
| Original Linear Regression | 4482.21 | 4618.91 | 0.9997 |
| Original XGBoost | 3843.48 | 4141.59 | 0.9999 |
| Original Random Forest | 3690.29 | 4018.50 | 0.9672 |
| Improved Linear Regression | 40.81 | 70.30 | 0.9997 |
| Improved XGBoost | 24.77 | 35.82 | 0.9999 |
| Improved Random Forest | 516.77 | 748.47 | 0.9672 |
| LSTM Neural Network | 15.23 | 22.45 | 0.9999 |

## 4.2 Key Findings

1. **Model Performance**

   - LSTM Neural Network outperforms all other models with lowest MAE and RMSE
   - Original Random Forest shows better performance than original XGBoost and Linear Regression
   - Improved models show significant performance enhancement over original models
   - Linear Regression shows the highest error rates among original models

2. **Clustering Analysis**

   - Moderate cluster separation (Silhouette Score: 0.413)
   - Clear distinction in demand patterns across clusters
   - Cluster 0 shows highest demand and temperature
   - Cluster 2 shows lowest demand and temperature
   - Cluster 3 shows highest humidity levels
   - Temporal patterns (hour) show similar distribution across clusters

3. **LSTM Advantages**

   - Better at capturing long-term dependencies in time series data
   - Superior performance in handling sequential patterns
   - More robust to noise and outliers
   - Can learn complex non-linear relationships
   - Particularly effective for electricity demand forecasting due to its ability to capture temporal patterns

4. **Feature Importance**

   - Weather conditions (especially temperature and humidity) are strong predictors
   - Temporal features (hour, day of week) show significant impact
   - City-specific patterns are important for accurate predictions

5. **Anomaly Detection**

   - Multiple methods (Z-score, IQR, Isolation Forest) provide robust anomaly detection
   - Seasonal patterns affect anomaly thresholds
   - Weather events correlate with demand anomalies

# 5. Conclusion and Recommendations

## 5.1 Key Takeaways

1. LSTM Neural Network provides the most accurate predictions with lowest error rates
2. Weather and temporal features are crucial for forecasting
3. Multiple anomaly detection methods provide robust validation
4. Deep learning approaches show superior performance for time series forecasting

## 5.2 Recommendations

1. Use LSTM as the primary forecasting model
2. Consider ensemble methods combining LSTM with XGBoost
3. Regular model retraining with new data
4. Consider city-specific models for better accuracy
5. Implement real-time LSTM predictions for dynamic forecasting

## 5.3 Future Work

1. Implement deep learning architectures
2. Explore more sophisticated ensemble methods
3. Incorporate additional external factors
4. Develop real-time anomaly detection system

# 6. References

- Scikit-learn documentation
- XGBoost documentation
- TensorFlow/Keras documentation
- Pandas documentation