## ⌄  XGBoost for Prediciting Energy Usage in Buildings

### Notebook Overview

This notebook demonstrates:

1. **Data Preprocessing**:

   ○ Cleaning, aligning, and preparing building-specific data.

2. **Model Evaluation**:

   ○ Using a trained XGBoost model to predict energy consumption.

3. **Visualization**:

   ○ Comparing actual and predicted values with interactive plots.

**Use cases**:

- **Model Validation**: Ensure predictions align with ground truth data.
- **Error Analysis**: Identify where and why predictions deviate from reality.
- **Presentation**: Generate clear visuals for stakeholder communication.

```
 1 import pandas as pd
 2 import numpy as np
 3 import os
 4 import matplotlib.pyplot as plt
 5 import re
 6 import json
 7 import gc
 8 import psutil
 9 import xgboost as xgb
10 import random
11 import joblib
12 import plotly.express as px
13
14 from multiprocessing import Pool, cpu_count
15 from concurrent.futures import ProcessPoolExecutor, as_completed
16 from functools import partial
17
18 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
19 from sklearn.model_selection import RandomizedSearchCV
```

```
1 # if you are using Google Colab, you can easily mount the drive and access
2 # the data.
3 from google.colab import drive
4 drive.mount('/content/drive', force_remount=True)
```

    ⇉  Mounted at /content/drive

## ⌄  Load the datasets

```
1 # In order for you to use the data, you need to update these paths.
2 PROCESSED = f'/content/drive/MyDrive/Team-Fermata-Energy/processed_data/' # https://drive.google.com/drive/folders,
3 BUILDINGS = f'{PROCESSED}processed_weather_load_w_timestamp/' # https://drive.google.com/drive/folders/1kW3Ip5_xm6I
```

```
1 with open(f'{PROCESSED}subset20_data.json', 'r') as test_train_file:
2     test_train_ids = json.load(test_train_file)
3
4 train_ids = [int(bldg_id.replace('.csv', '')) for bldg_id in test_train_ids['train_bldg_ids']]
5 test_ids = [int(bldg_id.replace('.csv', '')) for bldg_id in test_train_ids['test_bldg_ids']]
```

```
1 df_metadata = pd.read_csv(f"{PROCESSED}md_one_hot_encoded_subset20.csv")
2 df_metadata.head(20)
```

| | bldg_id | in.state | in.vintage | in.sqft | in.building_america_climate_zone_Cold | in.building_america_climate_zo |
|---|---|---|---|---|---|---|
| 0 | 105885 | 10 | 3 | 750000.0 | 0 | |
| 1 | 305819 | 40 | 2 | 150000.0 | 0 | |
| 2 | 305934 | 40 | 4 | 350000.0 | 0 | |
| 3 | 317044 | 40 | 3 | 350000.0 | 0 | |
| 4 | 32 | 1 | 6 | 37500.0 | 0 | |
| 5 | 64 | 1 | 0 | 37500.0 | 0 | |
| 6 | 103 | 1 | 4 | 75000.0 | 0 | |
| 7 | 112 | 1 | 3 | 7500.0 | 0 | |
| 8 | 277 | 1 | 0 | 37500.0 | 0 | |
| 9 | 355 | 1 | 7 | 17500.0 | 0 | |
| 10 | 363 | 1 | 5 | 37500.0 | 0 | |
| 11 | 379 | 1 | 1 | 17500.0 | 0 | |
| 12 | 417 | 1 | 0 | 17500.0 | 0 | |
| 13 | 530 | 1 | 2 | 17500.0 | 0 | |
| 14 | 575 | 1 | 3 | 3000.0 | 0 | |
| 15 | 611 | 1 | 5 | 7500.0 | 0 | |
| 16 | 633 | 1 | 6 | 37500.0 | 0 | |
| 17 | 864 | 1 | 4 | 37500.0 | 0 | |
| 18 | 1025 | 1 | 4 | 75000.0 | 0 | |
| 19 | 1215 | 1 | 3 | 3000.0 | 0 | |

20 rows × 40 columns

## If you want to load in the model and just evaluate, you can do so here!

```
1 model = joblib.load(f'{PROCESSED}xgb_model2.pkl')
```

## Functions for Training the Model

### Training the Model

This function prepares data and trains the XGBoost model for energy consumption prediction.

**Steps**:

1. Train the model: Set hyperparameters, use early stopping, and monitor performance with metrics like RMSE or MAE.

2. Analyze feature importance: Understand key drivers of predictions.

**Outputs**:

- A trained model

```python
1 def preprocess_bldg_optimized(bldg_id, df_metadata):
2     """
3     Preprocesses data for a single building.
4
5     Parameters:
6     - bldg_id (int): Building ID.
7     - df_metadata (DataFrame): Metadata DataFrame.
8
9     Returns:
10    - X (DataFrame): Feature DataFrame.
11    - y (Series): Target variable.
12    """
13    try:
14        # Load CSV with optimized data types
15        df_bldg = pd.read_csv(
16            f"{BUILDINGS}{bldg_id}.csv",
17            dtype={
18                'bldg_id': 'int32',
19                'minute': 'int8',
20                'out_electricity_total_energy_consumption': 'float32',
21                # Add other columns with appropriate types if known
22                # Example:
23                # 'temperature': 'float32',
24                # 'humidity': 'float32',
25                # 'heat_index': 'float32',
26                # 'location': 'category',
27            }
28        )
29
30        # Clean column names
31        df_bldg.columns = [re.sub(r"[^A-Za-z0-9_]+", "_", col) for col in df_bldg.columns]
32
33        # Filter rows where 'minute' == 0
34        df_bldg = df_bldg[df_bldg['minute'] == 0]
35
36        # Merge with metadata
37        bldg_metadata = df_metadata[df_metadata['bldg_id'] == bldg_id]
38        df_bldg = df_bldg.merge(bldg_metadata, on='bldg_id', how='left')
39
40        # Prepare features and target
41        y = df_bldg['out_electricity_total_energy_consumption']
42        X = df_bldg.drop(columns=['out_electricity_total_energy_consumption', 'timestamp', 'bldg_id'])
43
44        return X, y
45    except Exception as e:
46        print(f"Error processing building ID {bldg_id}: {e}")
47        return pd.DataFrame(), pd.Series()
48
49 def train_in_chunks_optimized(df_metadata, train_bldg_ids, best_param):
50     """
51     Trains an XGBoost model on data from multiple buildings using GPU acceleration.
52
53     Parameters:
54     - df_metadata (DataFrame): Metadata DataFrame.
55     - train_bldg_ids (list): List of building IDs for training.
56     - best_param (dict): Best hyperparameters for XGBoost.
57
```

```python
58      Returns:
59      - model (XGBRegressor): Trained XGBoost model.
60      """
61      # Determine number of parallel processes
62      n_jobs = max(cpu_count() - 1, 1)  # Reserve one core for the system
63      print(f"Using {n_jobs} parallel processes for data preprocessing.")
64
65      # Initialize multiprocessing Pool
66      with Pool(processes=n_jobs) as pool:
67          # Partial function to pass df_metadata
68          func = partial(preprocess_bldg_optimized, df_metadata=df_metadata)
69
70          # Map the preprocessing function to building IDs
71          results = pool.map(func, train_bldg_ids)
72
73      # Filter out any failed preprocessing results
74      results = [res for res in results if not res[0].empty]
75
76      if not results:
77          raise ValueError("No data was successfully preprocessed.")
78
79      # Concatenate all preprocessed data
80      X_all, y_all = zip(*results)
81      X_all = pd.concat(X_all, ignore_index=True)
82      y_all = pd.concat(y_all, ignore_index=True)
83
84      # Clean up intermediate results
85      del results
86      gc.collect()
87
88      print("Starting model training on GPU...")
89
90      # Initialize and train the XGBoost model with GPU support
91      model = xgb.XGBRegressor(
92          tree_method='hist',           # Use 'hist' for faster training with GPU via 'device'
93          device='cuda',                # Specify to use GPU
94          n_jobs=-1,                    # Utilize all available CPU cores for data preprocessing
95          enable_categorical=True,      # Enable categorical feature support
96          **best_param,                 # Additional hyperparameters
97          reg_alpha=1.0,
98          reg_lambda=1.0,
99          random_state=42,
100         verbosity=1                   # Set to 1 for basic logging
101     )
102
103     # Fit the model
104     model.fit(X_all, y_all, verbose=True)
105
106     # Monitor memory usage
107     process = psutil.Process()
108     memory_usage = process.memory_info().rss
109     print(f"Memory Usage After Training: {memory_usage / (1024 ** 2):.2f} MB")
110
111     # Clean up to free memory
112     del X_all
113     del y_all
114     gc.collect()
115
116     return model
117
```

```python
1 model = xgb.XGBRegressor(
2     tree_method='hist',           # Use 'hist' for faster training with GPU via 'device'
3     device='cuda',                # Specify to use GPU
4     n_jobs=-1,                    # Utilize all available CPU cores for data preprocessing
```

```
5     enable_categorical=True,      # Enable categorical feature support
6     reg_alpha=1.0,
7     reg_lambda=1.0,
8     random_state=42,
9     verbosity=1                   # Set to 1 for basic logging
10 )
```

```
1 best_param = {'subsample': 0.8, 'n_estimators': 300, 'max_depth': 6, 'learning_rate': 0.01, 'colsample_bytree': 0.8
2 model = train_in_chunks_optimized(df_metadata, train_ids, best_param)
```

> Using 11 parallel processes for data preprocessing.
> Starting model training on GPU...
> Memory Usage After Training: 37627.39 MB

## > Find best hyperparameters for the model.

I've already run this code so you don't have to.

```
[ ] ↳ 2 cells hidden
```

## ∨ Save and Test Model

### Evaluating the Model

This function tests the trained model on unseen data and calculates performance metrics.

**Steps**:

1. Load and preprocess test data: Align features with the trained model.
2. Predict and compare: Generate predictions and calculate metrics like SMAPE.
3. Monitor performance: Assess memory usage and model accuracy.

**Outputs**:

- Evaluation metrics and insights into prediction accuracy.

```
1 # Save the model
2 # joblib.dump(model, PATHGOESHERE)
```

> ['/content/drive/MyDrive/Team-Fermata-Energy/processed_data/xgb_model2.pkl']

```
1 df_test_bldg = pd.read_csv(f"{BUILDINGS}32.csv")
2 df_test_bldg.columns = [re.sub(r"[^A-Za-z0-9_]+", "_", col) for col in df_test_bldg.columns]
3 df_test_bldg.columns
```

> Index(['timestamp', 'out_electricity_total_energy_consumption',
>        'Dry_Bulb_Temperature_C_', 'Relative_Humidity_', 'heat_index', 'minute',
>        'hour', 'day', 'month', 'is_weekday', 'is_holiday', 'max_load_hourly',
>        'min_load_hourly', 'max_temp_hourly', 'min_temp_hourly', 'bldg_id'],
>       dtype='object')

```
1 # Pre-create numpy arrays to store the metrics
2 # num_test_ids = len(test_ids)  # Replace test_ids with your list of test building IDs
3 smape_values = []
4
5 def calculate_smape(y_true, y_pred):
6     """
7     Calculate Symmetric Mean Absolute Percentage Error (SMAPE).
8
```

```python
 9      Parameters:
10          y_true: Actual values.
11          y_pred: Predicted values.
12
13      Returns:
14          SMAPE value as a percentage.
15      """
16      denominator = (np.abs(y_true) + np.abs(y_pred)) / 2
17      diff = np.abs(y_true - y_pred)
18      smape = np.mean(diff / denominator) * 100  # Percentage
19      return smape
20
21 def evaluate_model(model, df_metadata, test_ids):
22      """
23      Loops through each test building ID, evaluates the model, and stores metrics.
24
25      Parameters:
26          model: Trained XGBoost model.
27          df_metadata: DataFrame containing building metadata.
28          test_ids: List of test building IDs.
29      """
30      for idx, bldg_id in enumerate(test_ids):
31          # Load and preprocess test building data
32          df_test_bldg = pd.read_csv(f"{BUILDINGS}{bldg_id}.csv")
33          df_test_bldg.columns = [re.sub(r"[^A-Za-z0-9_]+", "_", col) for col in df_test_bldg.columns]
34
35          df_test_bldg = df_test_bldg[df_test_bldg['minute'] == 0]
36
37          # Merge metadata
38          test_bldg_metadata = df_metadata[df_metadata['bldg_id'] == bldg_id]
39          df_test_bldg = df_test_bldg.merge(test_bldg_metadata, on='bldg_id', how='left')
40
41          # Prepare features (X_test) and target (y_test)
42          y_test = df_test_bldg['out_electricity_total_energy_consumption']
43          X_test = df_test_bldg.drop(columns=['out_electricity_total_energy_consumption', 'timestamp', 'bldg_id'])
44
45          # print(y_test.describe())
46          # Predict and evaluate metrics
47          pred = model.predict(X_test)
48          smape = calculate_smape(y_test, pred)
49          smape_values.append(smape)
50
51          # Monitor memory usage
52          process = psutil.Process()
53          memory_usage = process.memory_info().rss
54          print(f"Memory Usage: {memory_usage / (1024 ** 2):.2f} MB")
55
56          # Clean up memory
57          del df_test_bldg, X_test, y_test, pred, process, memory_usage
58          gc.collect()
```

```python
1 # Example usage
2 random.seed(521)
3 evaluate_model(model, df_metadata, random.choices(test_ids, k = 5))
4
5 # performance metrics
6 smape_array = np.array(smape_values)
7
8 print(f'Mean SMAPE: {mean_smape}')
```

```
Memory Usage: 19700.38 MB
Memory Usage: 19700.38 MB
Memory Usage: 19700.38 MB
Memory Usage: 19700.38 MB
Memory Usage: 19700.38 MB
```

```
      Mean SMAPE: 21.48865254136774
```

## ∨ Visualizations and Misc

### Align Features for Model Compatibility

The `align_features` function ensures that the test data matches the trained model's expected input format. This is critical for avoiding feature name mismatches, which occur when:

- The test data contains columns that were not part of the training data.
- The test data is missing columns present during model training.

**Key steps**:

1. **Add Missing Columns**: Columns that are in the model's feature list but not in the test data are added with default values (e.g., `0` for one-hot encoded features).
2. **Drop Extra Columns**: Any columns in the test data but not required by the model are removed.
3. **Order Matching**: Ensures that the columns are in the same order as the model's training data.

This function ensures smooth prediction and prevents runtime errors due to feature mismatch.

```python
 1 def align_features(X_test, model):
 2     """
 3     Aligns test data features with the features expected by the model.
 4
 5     Parameters:
 6         X_test (pd.DataFrame): Test data features.
 7         model: Trained model (XGBoost or similar).
 8
 9     Returns:
10         pd.DataFrame: Aligned test data.
11     """
12     # Get the feature names from the model
13     model_features = model.get_booster().feature_names
14
15     # Add missing columns to X_test
16     for col in model_features:
17         if col not in X_test.columns:
18             X_test[col] = 0  # Default value for missing features
19
20     # Drop extra columns not in the model
21     X_test = X_test[model_features]
22
23     return X_test
24
25
26 def visualize_time_series(df_metadata, bldg_id, model):
27     """
28     Load and visualize time series data for energy consumption and model predictions.
29
30     Parameters:
31         df_metadata (pd.DataFrame): DataFrame containing building metadata.
32         bldg_id (str): Building ID for which the data is visualized.
33         model: Trained XGBoost model for predictions.
34     """
35     # Load the specific building data
36     file_path = f"{BUILDINGS}{bldg_id}.csv"
37     df_bldg = pd.read_csv(file_path)
38
39     # Clean column names
```

```
40      df_bldg.columns = [re.sub(r"[^A-Za-z0-9_]+", "_", col) for col in df_bldg.columns]
41      print(df_bldg.columns)
42
43      # Convert timestamp to datetime
44      df_bldg['timestamp'] = pd.to_datetime(df_bldg['timestamp'])
45
46      # Filter rows where minute == 0
47      df_bldg = df_bldg[df_bldg['minute'] == 0]
48
49      # Prepare features for prediction
50      X_test = df_bldg.drop(columns=['out_electricity_total_energy_consumption', 'timestamp', 'bldg_id'])
51      X_test = align_features(X_test, model)
52
53      # Predict using the model
54      df_bldg['Predicted_Energy_Consumption'] = model.predict(X_test)
55
56      # Melt the DataFrame for easier plotting
57      df_long = df_bldg.melt(
58          id_vars='timestamp',
59          value_vars=[
60              'out_electricity_total_energy_consumption',
61              'Predicted_Energy_Consumption'
62          ],
63          var_name='Measurement',
64          value_name='Value'
65      )
66
67      # Replace specific measurement names for better legend labels
68      df_long['Measurement'] = df_long['Measurement'].replace({
69          'out_electricity_total_energy_consumption': 'Actual Energy Consumption',
70          'Predicted_Energy_Consumption': 'Predicted Energy Consumption'
71      })
72
73      # Create the time series plot
74      fig = px.line(
75          df_long,
76          x='timestamp',
77          y='Value',
78          color='Measurement',
79          labels={'Value': 'Measurement Value', 'timestamp': 'Time'},
80          title=f"Time Series Data for Building ID: {bldg_id} (Actual vs Predicted)"
81      )
82
83      # Show the plot
84      fig.show()
```

## ∨  Visualize Time Series Data (Actual vs Predicted)

The `visualize_time_series` function provides a comprehensive view of the model's performance by comparing actual energy consumption against the predicted values.

**Key components**:

1. **Data Preparation**:

   - Load the building-specific dataset and filter rows where `minute == 0` for consistent granularity.
   - Align features with the trained model using the `align_features` function.

2. **Predictions**:

   - The model predicts energy consumption using preprocessed test data.
   - Predicted values are added as a new column to the dataset.
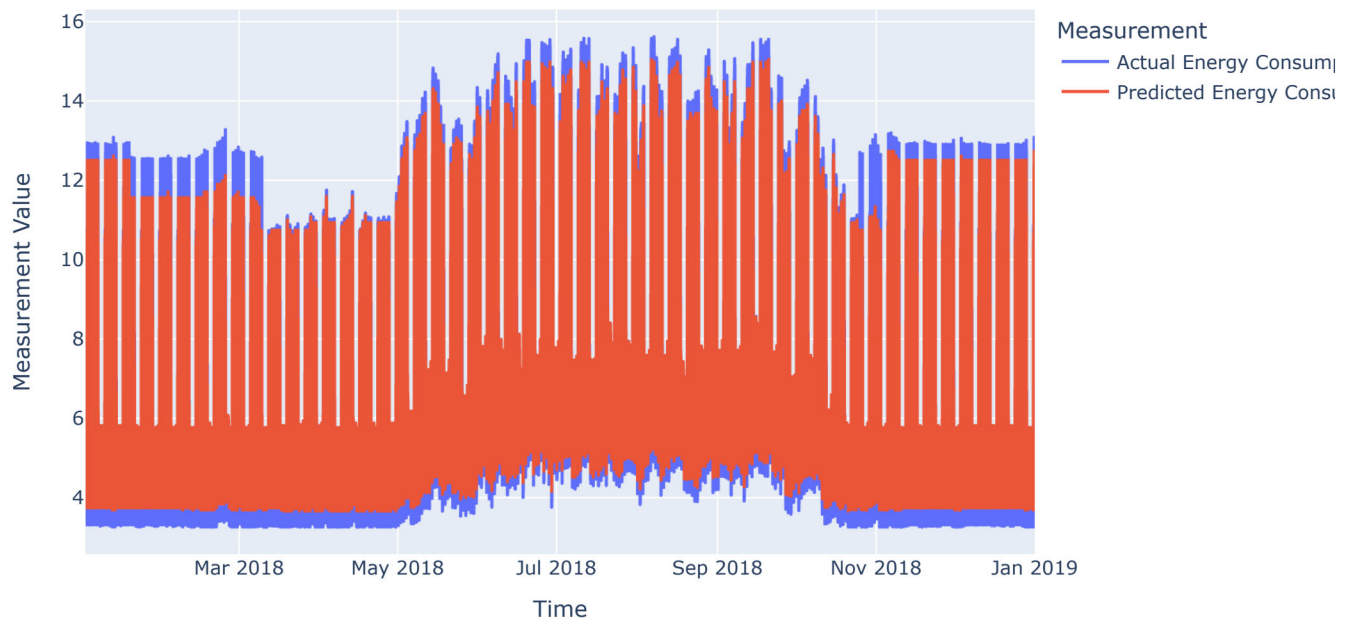
3. **Visualization**:

- A line plot compares actual and predicted energy consumption over time.
- Interactive legends allow users to focus on specific measurements.

This visualization helps identify patterns, trends, and areas where the model's predictions deviate from the actual values.

```
1 visualize_time_series(df_metadata, '32', model)
```

```
Index(['timestamp', 'out_electricity_total_energy_consumption',
       'Dry_Bulb_Temperature_C_', 'Relative_Humidity_', 'heat_index', 'minute',
       'hour', 'day', 'month', 'is_weekday', 'is_holiday', 'max_load_hourly',
       'min_load_hourly', 'max_temp_hourly', 'min_temp_hourly', 'bldg_id'],
      dtype='object')
```
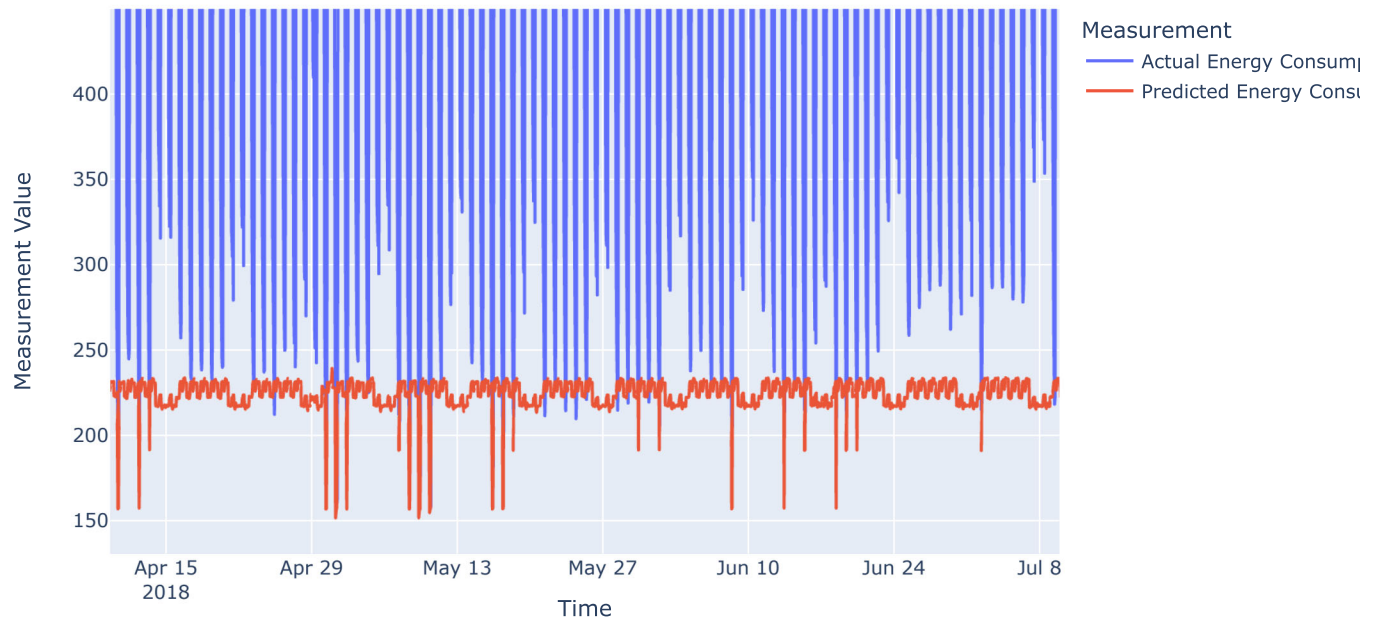
### Time Series Data for Building ID: 32 (Actual vs Predicted)



```
1 visualize_time_series(df_metadata, '105885', model)
```

```
Index(['timestamp', 'out_electricity_total_energy_consumption',
       'Dry_Bulb_Temperature_C_', 'Relative_Humidity_', 'heat_index', 'minute',
       'hour', 'day', 'month', 'is_weekday', 'is_holiday', 'max_load_hourly',
       'min_load_hourly', 'max_temp_hourly', 'min_temp_hourly', 'bldg_id'],
      dtype='object')
```

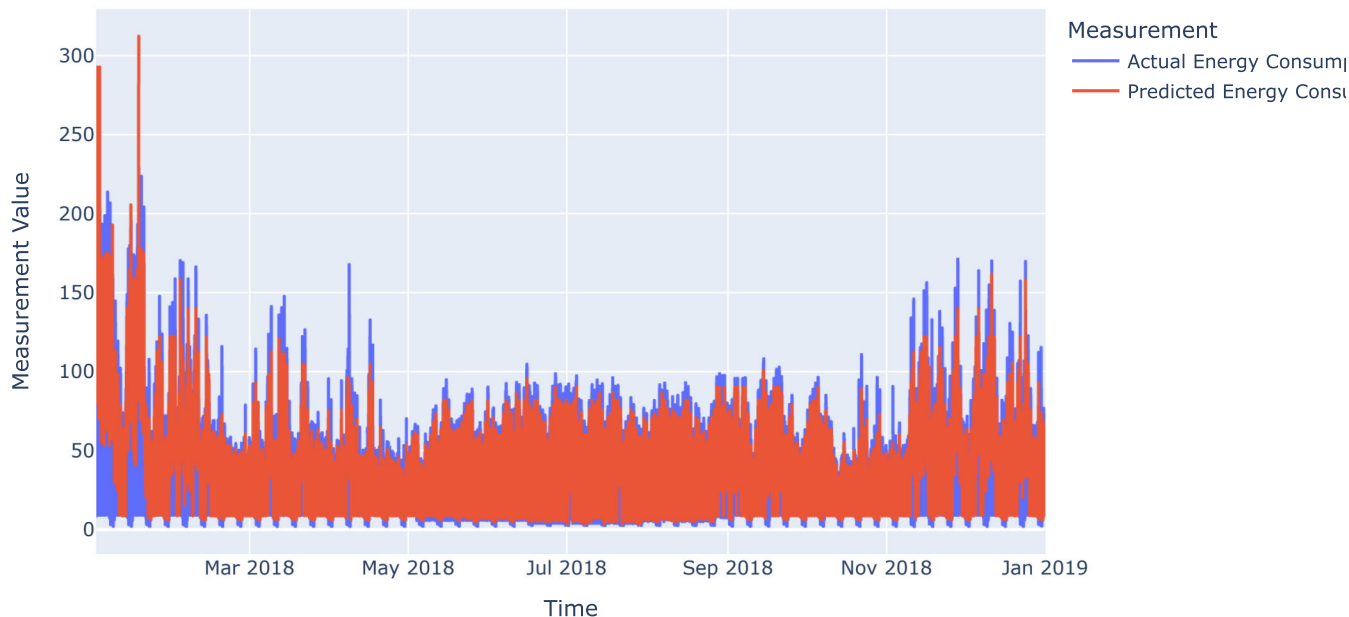Time Series Data for Building ID: 105885 (Actual vs Predicted)



```
1 visualize_time_series(df_metadata, '1025', model)
```

```
Index(['timestamp', 'out_electricity_total_energy_consumption',
       'Dry_Bulb_Temperature_C_', 'Relative_Humidity_', 'heat_index', 'minute',
       'hour', 'day', 'month', 'is_weekday', 'is_holiday', 'max_load_hourly',
       'min_load_hourly', 'max_temp_hourly', 'min_temp_hourly', 'bldg_id'],
      dtype='object')
```

### Time Series Data for Building ID: 1025 (Actual vs Predicted)



```
1 feature_importances = model.feature_importances_
2 features = []
3 features = model.get_booster().feature_names
4
5 # Create a DataFrame for better handling of the data
6 data = pd.DataFrame({
7     'Feature': features,
8     'Importance': feature_importances
9 })
10
11 # Sort features by importance for better visualization
12 data = data.sort_values(by='Importance', ascending=True)
13
14 # Plot the feature importances using Plotly
15 fig = px.bar(data, x='Importance', y='Feature', orientation='h',
16              title='Feature Importances for Load Forecasting',
17              labels={'Importance': 'Feature Importance', 'Feature': 'Feature'}
18              text='Importance')
19
20 # Improve layout for better readability
21 fig.update_layout(
22     xaxis_title="Feature Importance",
23     yaxis_title="Feature",
24     title_x=0.5,  # Center the title
25     font=dict(size=12),
26     showlegend=False,
27     margin=dict(l=150, r=20, t=50, b=50),  # Adjust left margin for long label
28     height=400 + 20 * len(features)  # Dynamically adjust height for label siz
29 )
30
31 # Add better formatting for the text
```

```
32 fig.update_traces(texttemplate='%{text:.2f}', textposition='outside')
33
34 # Display the plot
35 fig.show()
```

## Feature Importances for Load Forecasting