**L06: IIoT Time Series Forecasting Lab Report**

## 1. Introduction

This lab explores the application of time-series forecasting techniques to real-world Industrial Internet of Things (IIoT) temperature data. Time series forecasting plays a critical role in IIoT environments, enabling predictive maintenance, resource optimization, and anomaly detection. The primary objectives of this lab were to:

- Preprocess IIoT temperature data for forecasting

- Engineer meaningful features that capture temporal patterns

- Select and train appropriate forecasting models

- Evaluate model performance using multiple metrics

- Apply generative modeling techniques to enhance the dataset

The lab demonstrates a complete machine learning pipeline for IIoT time series data, from raw data to performance evaluation, with a focus on enhancing prediction accuracy through feature engineering and data augmentation using generative models.

## 2. Data Preparation

### Dataset Overview

The dataset contained 97,606 temperature records from IoT sensors with the following structure:

- id: Unique identifier for each record

- room_id/id: Identifier for the room or location

- noted_date: Timestamp of the reading (in DD-MM-YYYY HH:MM format)

- temp: Temperature reading

- out/in: Indicator of whether the reading was from inside or outside

### Preprocessing Steps

Several preprocessing steps were performed to prepare the data for forecasting:

1. **Column Renaming**: The noted_date and temp columns were renamed to timestamp and temperature for clarity and consistency.

2. **Missing Value Handling**: The dataset was checked for missing values. Although no missing values were found in the original dataset, forward and backward filling strategies were implemented to ensure robustness.

3. **Timestamp Conversion**: The timestamp column was converted to datetime format and set as the index to facilitate time-based operations.

4. **Outlier Detection and Handling**: The Interquartile Range (IQR) method was used to identify outliers in the temperature data. Rather than removing these outliers, which would disrupt the time series continuity, they were capped at the lower and upper bounds (Q1 - 1.5*IQR and Q3 + 1.5*IQR, respectively).

5. **Train-Test Split**: The data was split into training (80%) and testing (20%) sets, respecting the temporal nature of the data by maintaining chronological order.

**Visualizations**

The temperature time series visualization revealed both short-term fluctuations and longer-term patterns, confirming the suitability of time series forecasting approaches. The train-test split visualization demonstrated that the testing period contained similar patterns to the training period, suggesting that a well-trained model should generalize effectively.

**3. Feature Engineering**

Feature engineering is crucial for improving forecasting accuracy, especially when dealing with time series data that contains multiple cyclical and trend components. Four custom features were created based on domain knowledge of IoT temperature patterns:

**Custom Features**

1. **Hour of Day**: This feature captures diurnal temperature patterns, as temperatures typically follow a daily cycle with highs during daylight hours and lows at night. By extracting the hour (0-23) from the timestamp, the model can learn these regular daily fluctuations.

2. **Day of Week**: This feature identifies weekly patterns that might exist due to different usage patterns of the monitored spaces across weekdays and weekends. The day of week (0-6, where 0 is Monday) was extracted from the timestamp.

3. **Previous Temperature (lag-1)**: This feature incorporates temporal dependency by using the previous temperature reading to predict the next value. Temperature values show strong autocorrelation, making this a highly predictive feature.

4. **Moving Average (24-hour window)**: This feature captures longer-term trends by calculating a rolling 24-hour moving average. It helps smooth out short-term fluctuations and highlights the underlying trend.

**Feature Significance Analysis**

Correlation analysis and feature importance evaluation revealed that the previous temperature and moving average features were the most predictive, with correlation coefficients of 0.93 and 0.89 with the target temperature, respectively. This aligns with domain knowledge that recent temperature history is often the strongest predictor of future temperature.

The hour of day feature showed moderate correlation (0.41), confirming the presence of daily patterns. The day of week feature showed the weakest correlation (0.12), suggesting weekly patterns were less prominent in this dataset.

## 4. Model Selection and Training

**Model Selection**

For this lab, a Random Forest Regressor was chosen for time series forecasting due to several advantages:

1. **Robustness to outliers**: Despite our outlier handling, Random Forest's inherent robustness provides an additional layer of protection against extreme values.

2. **Ability to capture non-linear relationships**: Temperature patterns can involve complex non-linear relationships that Random Forest can effectively model.

3. **Feature importance information**: Random Forest provides valuable insights into which features contribute most to prediction accuracy.

4. **Minimal hyperparameter tuning**: Random Forest performs well with default parameters, making it suitable for initial modeling efforts.

While Nixtla's AutoML was originally considered for model selection, access issues led to the adoption of the Random Forest approach, which proved to be effective for this task.

**Training Process**

The model was trained using the following steps:

1. Features and target were prepared from the training data

2. Features were standardized using StandardScaler to ensure all features contributed equally

3. A Random Forest model with 100 estimators was trained on the scaled features

4. The model's feature importance was analyzed to understand the relative contribution of each feature

## 5. Evaluation

**Performance Metrics**

The model was evaluated using multiple metrics to provide a comprehensive assessment of its forecasting performance:

| Metric | Original Model | Description |
|--------|----------------|-------------|
| MAE | 0.7314 | Average absolute difference between predicted and actual temperatures |
| MSE | 0.9741 | Average squared difference, penalizing larger errors more heavily |
| RMSE | 0.9870 | Square root of MSE, in the same units as the temperature |
| MAPE | 2.45% | Percentage error relative to the actual temperature |
| MASE | 0.2939 | Ratio of MAE to the MAE of a naive forecast |

The MASE value of 0.2939 is particularly notable as it indicates the model significantly outperformed a naive forecast (using the previous temperature as the prediction). A MASE < 1 indicates superior performance compared to the naive approach, and our value of 0.2939 suggests the model is approximately 3.4 times more accurate.

The MAPE of 2.45% demonstrates high prediction accuracy relative to the actual temperature values, indicating that on average, predictions were within about 2.5% of the true temperatures.

**Cross-Validation Results**

To ensure the model's robustness across different time periods, time series cross-validation was performed using a rolling-origin approach. This involved creating multiple train-validation splits that moved forward in time, training a model on each split, and evaluating it on the validation set.

The average cross-validation MAE was 1.4478, with a standard deviation of 0.0521 across folds. The higher MAE in cross-validation compared to the test set evaluation suggests

some variability in the model's performance across different time periods, which is common in time series forecasting.

**Residual Analysis**

Residual analysis showed that prediction errors were approximately normally distributed around zero, with no obvious patterns that would indicate systematic bias in the model. This confirms that the model captured the underlying patterns in the data effectively.

## 6. Generative Modeling

**VAE Architecture and Implementation**

A Variational Autoencoder (VAE) was implemented to generate synthetic temperature data. VAEs were chosen over GANs due to their computational efficiency and ability to learn a smooth latent space representation of the data.

The VAE architecture consisted of:

- An encoder network with two hidden layers (32 and 16 neurons) that compressed the temperature data into a 4-dimensional latent space

- A sampling mechanism that added variability through the reparameterization trick

- A decoder network with two hidden layers (16 and 32 neurons) that reconstructed temperature data from the latent representation

The model was trained for 20 epochs using a custom loss function that combined reconstruction loss (mean squared error) and KL divergence to ensure a structured latent space.

**Synthetic Data Generation and Integration**

After training, the VAE was used to generate 1,000 synthetic temperature records by sampling from the latent space and passing the samples through the decoder. The synthetic timestamps were created by continuing from the end of the training set with hourly intervals.

The synthetic data preserved the same feature structure as the original data, with hour of day, day of week, previous temperature, and moving average features calculated for each synthetic record. This ensured compatibility with the forecasting model.

The visual comparison between original and synthetic data showed that the VAE successfully captured the general distribution characteristics of the temperature data, with similar ranges and variability.

**Performance Improvement Analysis**

The original training set was augmented with the synthetic data, and the forecasting model was retrained on this combined dataset. The performance comparison showed:

| Metric | Original Model | Augmented Model | Improvement |
|--------|----------------|-----------------|-------------|
| MAE | 0.7314 | 0.7294 | 0.28% |
| MSE | 0.9741 | 0.9715 | 0.27% |
| RMSE | 0.9870 | 0.9856 | 0.14% |
| MAPE | 2.45% | 2.44% | 0.41% |
| MASE | 0.2939 | 0.2931 | 0.27% |

The modest improvement of 0.28% in MAE suggests that the original dataset was already fairly representative, but the synthetic data still provided a small benefit. This aligns with expectations, as the primary benefit of synthetic data is often observed in scenarios with limited training data, while our dataset was relatively large (~78,000 training examples).

## 7. Individual Reflection

**Challenges and Adaptations**

Several challenges were encountered during this lab:

1. **API Access Issues**: The original plan to use Nixtla's AutoML for model selection had to be adapted due to API access issues. This led to the implementation of a Random Forest approach, which proved effective and provided valuable insights through feature importance analysis.

2. **Timestamp Format Handling**: The dataset used a non-standard date format (DD-MM-YYYY) which required careful handling during preprocessing to ensure correct temporal ordering.

3. **VAE Hyperparameter Tuning**: Finding the right balance between reconstruction loss and KL divergence in the VAE loss function required experimentation to generate realistic synthetic data.

These challenges provided valuable learning opportunities in adapting techniques to specific dataset characteristics and constraints.

**Learning Experience**

This lab provided hands-on experience with the complete lifecycle of a time series forecasting project, from data preparation to model evaluation and enhancement. Key learnings include:

1. The importance of feature engineering in time series forecasting, especially creating features that capture cyclical patterns and temporal dependencies

2. The value of multiple evaluation metrics to understand different aspects of model performance

3. The application of generative models for data augmentation in time series contexts

4. The need for rolling-origin cross-validation to ensure model robustness across different time periods

**Future Improvements**

Several potential improvements could enhance the forecasting performance:

1. **Additional Features**: Incorporating external data such as weather conditions could improve accuracy, especially for outdoor temperature readings.

2. **Advanced Models**: Exploring deep learning approaches like LSTMs or Transformer-based models could capture more complex temporal dependencies.

3. **Hierarchical Forecasting**: Creating separate models for indoor and outdoor readings (based on the 'out/in' column) might improve performance by accounting for different patterns in these environments.

4. **Hybrid Approach**: Combining statistical methods (e.g., ARIMA) with machine learning models could leverage the strengths of both approaches.

5. **More Sophisticated Generative Models**: Exploring conditional VAEs or time-series specific GANs could generate more realistic synthetic data with stronger temporal coherence.

**8. Conclusion**

This lab successfully demonstrated the application of time series forecasting techniques to IIoT temperature data. The Random Forest model achieved high accuracy with a MAPE of 2.45% and a MASE of 0.2939, indicating strong predictive performance relative to actual temperatures and naive forecasting baselines.

Feature engineering proved crucial, with temporal features like previous temperature and moving averages providing significant predictive power. The generative modeling approach

using a VAE demonstrated the potential of synthetic data for enhancing forecasting performance, albeit with modest improvements in this case.

The methodology and insights from this lab can be applied to various IIoT forecasting scenarios, from energy consumption prediction to equipment maintenance scheduling, highlighting the value of time series forecasting in industrial settings.

## 9. References

1. "Forecasting: Principles and Practice" by Rob J Hyndman and George Athanasopoulos

2. "Deep Learning for Time Series Forecasting" by Jason Brownlee

3. "Auto-Encoding Variational Bayes" by Kingma and Welling (2014)

4. "Time Series Analysis by State Space Methods" by J. Durbin and S.J. Koopman