

Electricity Price Prediction and Smoothing for AEMO

Tanmay, Somani (n11485094)

IFN704 Assessment 3

Executive Summary

This project aimed to design and evaluate machine learning models to predict electricity prices and minimize price and demand spikes in the NSW region of the National Electricity Market (NEM). Addressing price volatility is crucial for ensuring market stability and facilitating well-informed decisions in the energy sector. This project leveraged machine learning models, specifically a Random Forest Classifier and an RNN-LSTM hybrid model, to effectively forecast electricity prices and identify price spikes. A central focus of the project was to detect and normalize sudden price changes, in order to ensure more stable market behaviour. Furthermore, Gaussian smoothing techniques were utilized to mitigate the influence of extreme fluctuations. The project outcomes demonstrated that combining machine learning models with effective feature engineering and post-processing methods, like smoothing, can provide a robust solution to the challenges posed by high price volatility in electricity markets.

Introduction

Electricity price volatility is a persistent challenge in the National Electricity Market (NEM), particularly in regions like New South Wales (NSW). Sudden price changes, known as price spikes, can have significant implications for energy producers, consumers, and grid stability. Accurate price forecasting is critical to mitigating these fluctuations, ensuring market stability, and enabling stakeholders to make informed decisions.

This project aimed to develop a robust predictive model for electricity prices in the NSW region of the NEM, focusing on minimizing price and demand spikes. Machine learning models, specifically a Random Forest Classifier and an RNN-LSTM hybrid model, were employed to address the inherent complexities of electricity price forecasting, such as temporal dependencies, non-linearity, and the high volatility of prices.

The Random Forest Classifier was primarily used for feature selection and to classify instances of price spikes. In contrast, the RNN-LSTM hybrid model was designed to capture temporal relationships and predict future prices effectively. Moreover, Gaussian smoothing techniques were integrated into the process to reduce the impact of extreme price fluctuations, providing a more accurate and stable forecast.

The integration of machine learning techniques, feature engineering, and smoothing methods demonstrated the effectiveness of such models in addressing high price volatility. This report presents the data collection and preprocessing steps, the modelling methodologies used, the analysis of results, and a discussion of the challenges encountered and future improvements. The goal was to develop a

solution capable of providing stable and realistic price predictions, thereby supporting the decision-making processes of market participants.

Literature Review

1. Machine Learning Models in Price Forecasting (Chattopadhyay et al., 2020)

Machine learning models have become increasingly prevalent in electricity price forecasting, particularly following the deregulation of electricity markets. A study carried out by Chattopadhyay et al. (2020) highlights how price fluctuations often arise from rapid changes in demand, unexpected outages in generation units, and variability in fuel prices. To handle such complexities, a range of machine learning approaches, such as **linear regression models**, **support vector machines (SVM)**, **decision trees**, and **neural networks**, have been adopted.

Regression models provide a fundamental approach to forecasting, but their accuracy can be limited in scenarios involving nonlinear dependencies and high volatility. **LSTMs** (Long Short-Term Memory Networks) and other recurrent neural networks have shown considerable promise in capturing these dependencies because of their inherent ability to remember patterns over time.

2. Impact of External Factors on Price Volatility (Yan et al., 2019)

A study by Yan et al. (2019) focused on integrating multiple sources of data, such as temperature, renewable generation levels, and interconnector flows, to enhance the accuracy of price forecasting models. This research demonstrated that incorporating weather information significantly improved the ability to forecast peak demand periods, leading to better predictions of price volatility.

Incorporating **temperature data** allowed models to anticipate spikes in electricity consumption due to increased use of heating or cooling systems during extreme weather conditions. This aligns with our decision to integrate weather data, enabling our model to more accurately capture the factors driving electricity demand in Sydney.

3. Importance of Smoothing Techniques and Feature Engineering (Wang et al., 2022)

Wang et al. (2022) discussed various feature engineering techniques that could significantly impact model performance. One crucial aspect of their research was **Gaussian smoothing**, which can remove noise from data, especially in the presence of volatile price movements. The study demonstrated that **smoothing techniques** could enhance spike detection by minimizing short-term random fluctuations that do not reflect true market conditions. We leveraged Gaussian smoothing to reduce the impact of transient noise, improving our model's spike detection capabilities.

4. Class Imbalance in Electricity Price Spikes (Vargas et al., 2021)

Vargas et al. (2021) highlighted one of the key challenges in electricity price forecasting—handling the imbalance in the dataset when predicting rare events such as price spikes. Spikes are infrequent compared to non-spike periods, leading to challenges in effectively training models to recognize these events without overwhelming bias towards the more frequent, non-spike class. Vargas et al. proposed methods such as **oversampling** and **cost-sensitive learning** to counter this problem.

In our model, the **Random Forest classifier** struggled to achieve high recall for spikes due to the class imbalance, a problem also highlighted in the literature. Addressing this challenge effectively would likely require additional techniques such as **SMOTE (Synthetic Minority Over-sampling Technique)**.

5. Machine Learning and Electricity Markets – A Review (Zhao et al., 2020)

Zhao et al. (2020) provided a comprehensive review of machine learning applications in electricity market studies. This review shed light on how different **neural network architectures**, such as **Convolutional Neural Networks (CNNs)** and **Recurrent Neural Networks (RNNs)**, are applied to price forecasting, load prediction, and spike detection. CNNs are particularly beneficial for extracting spatial patterns, whereas RNNs are effective in sequential or time-series problems like price forecasting.

The insights gained from Zhao et al.'s review informed our model selection process. Specifically, we opted for a hybrid RNN + LSTM model, which combines the advantages of both—capturing the temporal dependencies effectively while leveraging RNN's strengths in sequence processing.

Methodology

The project involved several phases, including data collection, preprocessing, feature engineering, model training, and evaluation. This section will elaborate on each of these steps in detail, showcasing the approaches used to handle the complexities of electricity price prediction and spike detection.

Data Collection and Preprocessing

The datasets used in the project were obtained from the **Australian Energy Market Operator (AEMO)**, which included:

- **Historical electricity prices**
- **Demand values**
- **Rooftop solar generation data**
- **Other market parameters such as hourly temperatures for capital cities**

These datasets covered multiple aspects of the market and provided a comprehensive view for analysing electricity price dynamics. The data spanned the entire year of 2022 and included different regions: **Adelaide, Brisbane, Sydney, Canberra, Melbourne, and Hobart**. These regions were initially labelled with respective data counts: **Adelaide (23000), Brisbane (40913), Sydney (66214), Canberra (70351), Melbourne (86282), and Hobart (94250)**. After careful analysis, it was found that **New South Wales (NSW)** exhibited both price and demand spikes, making it the most suitable candidate for further modelling (Figure 1).



Figure 1 Demand and Price spikes by region

Data Cleaning and Integration

The datasets were merged and cleaned to ensure consistency across different data types and remove any anomalies. This included:

- **Removing irrelevant columns** such as **LASTCHANGED**, **REGIONID**, and **INTERVAL_DATETIME** to focus only on essential data that contributes to the prediction models.
- **Handling missing values** using **forward and backward fill techniques** to ensure no gaps in the dataset. The reason behind using this method was to make sure we don't have random values in the missing values.
- **Managing outliers** through **data clipping** and **transformation** to reduce their impact on model training and prevent skewed predictions.

To facilitate subsequent analyses, the data was filtered based on **region ID** to focus only on NSW. Once the NSW region was selected, multiple data files were utilized to create a new dataset containing the most relevant features. These included:

- **Prices by State 2022**
- **Actual Demand by State 2022**
- **Forecast Demand by State 2022**
- **Rooftop Solar Production 2022**
- **Hourly Temperatures Capital Cities 2022**

Feature Engineering

Feature engineering played a critical role in making the model effective for electricity price prediction. The initial features derived from the datasets included:

- **OPERATIONAL_DEMAND:** Represents actual demand, which is a crucial feature as price is heavily influenced by real-time demand.
- **POWER:** Total power generated, helping to understand the supply side of the equation.
- **Lagged Demand and Price Values:** Lagged features, such as **demand_lag_1**, **demand_lag_3**, **price_lag_1**, and **price_lag_3**, were created to capture temporal dependencies in price and demand patterns.
- **Day_of_Week** and **Hour_of_Day:** Categorical features representing the specific day and hour to capture the cyclical nature of electricity consumption patterns.
- **Temperature:** Weather plays an important role in electricity demand, especially due to heating and cooling loads.

The features were then **scaled** using the **MinMaxScaler** to normalize them within a suitable range, which is necessary for the RNN and LSTM models that are sensitive to the scale of the input data.

Price Spike Threshold

An initial threshold of **20%** was established to identify price spikes. If the electricity price increased by more than **20%** within a given interval, it was considered a spike. This threshold was chosen based on historical data analysis, which indicated that most significant price movements involved changes greater than **20%**. Setting this threshold allowed us to:

- Detect the maximum number of relevant price spikes while avoiding small, less impactful fluctuations.
- Focus on major market events that would require involvements or are of interest to energy market participants.
- Reduce noise in the data, making the modelling task more manageable.

This choice of threshold ensured that only meaningful price changes were flagged, providing a practical basis for predicting and managing price volatility. It also aligned with typical price response strategies used in market trading, where only substantial changes justify energy trading decisions.

Feature Selection and Importance Analysis

After preprocessing the data, a **feature importance analysis** was conducted using the **Random Forest model**. The goal was to identify which features had the most predictive power for detecting price spikes. The following features were ranked as the most important(Fig 2):

- **lagged_demand** (Importance: 0.2609)
- **OPERATIONAL_DEMAND** (Importance: 0.1901)
- **POWER** (Importance: 0.1702)

- demand_error (Importance: 0.1435)
- temperature (Importance: 0.1249)
- hour_of_day (Importance: 0.0564)
- day_of_week (Importance: 0.0539)

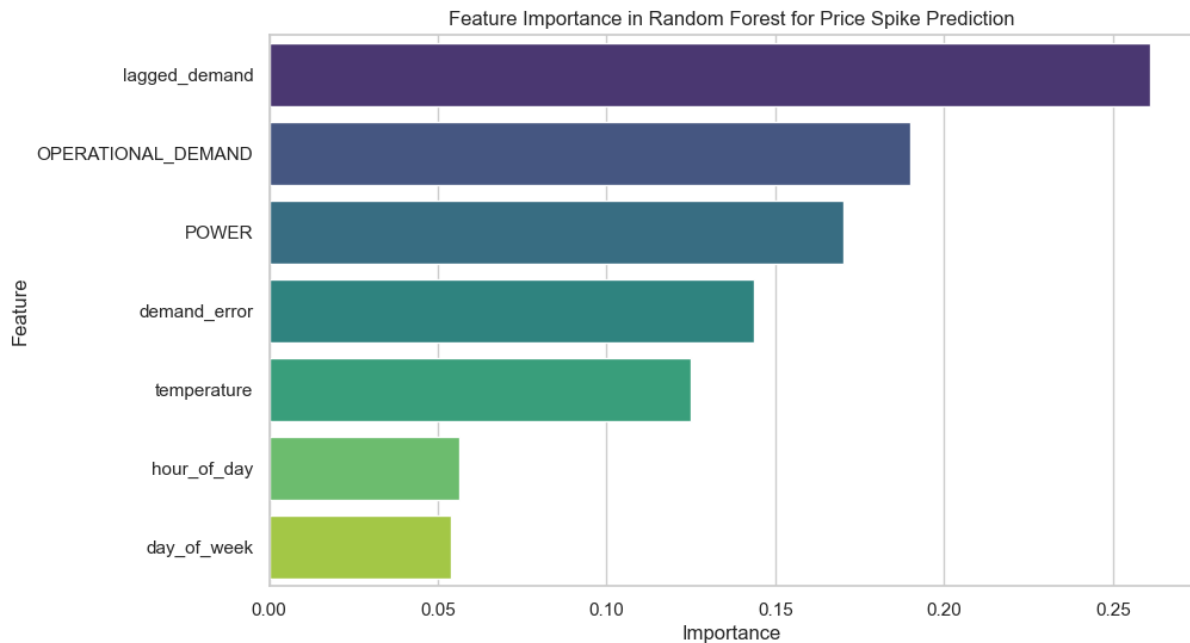


Figure 2 Feature Importance Rankings

These features provided insight into which factors were most influential in driving price changes and spikes. **Lagged demand** and **operational demand** were particularly important, as they captured immediate past behaviour and real-time dynamics, both of which are critical in forecasting electricity prices.

Model Development

The development of predictive models for electricity prices involved the implementation of two key models: the Random Forest Classifier and an RNN-LSTM hybrid model. Each of these models played a distinct role in addressing the unique challenges of electricity price forecasting and price spike detection.

The **Random Forest Classifier** was primarily used for binary classification of price spikes. The model was trained using some key features identified during feature importance analysis, such as **OPERATIONAL_DEMAND**, **POWER**, **lagged demand**, and **temperature**. By leveraging these features, the Random Forest provided valuable insights into the relationships between various factors and their contribution to price spikes. It was particularly effective in ranking the importance of different features, helping to identify the most significant contributors to price volatility. Although the Random Forest Classifier did not inherently account for temporal dependencies, it served as an essential base model for determining whether a price spike had occurred. This made it feasible for feature selection and initial spike identification.

The **RNN-LSTM hybrid model** was developed to handle the more complex problem of time-series forecasting of electricity prices. The architecture of the LSTM model consisted of several components. It included **two LSTM layers** designed to capture long-term dependencies in the sequential data. These layers were followed by a **SimpleRNN layer**, which was included to effectively model shorter-term sequential patterns that are often crucial for accurate forecasting. To reduce overfitting during training, a **Dropout layer** was included to randomly drop neurons. Finally, a **Dense layer** with a sigmoid activation function was used to produce the final output for binary classification of price spikes.

The RNN-LSTM model was trained on **80% of the dataset**, while the remaining **20%** was used for testing. The training process consisted of **20 epochs** with a **batch size of 64**, and a validation dataset was used to monitor the model's performance and prevent overfitting. The **binary cross-entropy loss function** was utilized, in combination with the **Adam optimizer**, to effectively minimize prediction errors during training. For effective time-series modelling, sequences of **10-time steps** were created using the `create_lstm_sequences` function. This sequence length was chosen to balance capturing relevant short-term trends while avoiding the introduction of excessive noise that might arise from longer sequences.

The performance of the RNN-LSTM model was evaluated using multiple metrics. **Accuracy** was used to assess the overall correctness of the model in classifying price spikes. The **Mean Absolute Error (MAE)** was used to quantify the average magnitude of the errors in the predictions, and the model achieved an MAE of approximately **4.15**, demonstrating reasonable prediction accuracy. Additionally, the **Symmetric Mean Absolute Percentage Error (SMAPE)** was used to provide insight into the percentage deviation of the predictions from actual values. Overall, the RNN-LSTM hybrid model demonstrated a strong ability to learn temporal dependencies, making it suitable for forecasting price variations effectively in the context of the NEM.

Training and Evaluation Process

Model training involved multiple iterations to find the best hyperparameters for each model. For the RNN-LSTM hybrid model, the **lookback period** was particularly challenging to determine, as it significantly impacted the model's ability to capture meaningful patterns without introducing unnecessary noise. **Hyperparameter tuning** was also conducted to balance performance with computational efficiency, including:

- **Batch size and number of epochs** for training.
- **Number of LSTM units and dropout layers** to prevent overfitting.

The **Random Forest model** was assessed primarily for its ability to classify price spikes. The **precision, recall, and F1-score** were calculated to determine the model's capability in detecting both spike and non-spike events accurately. The overall **accuracy** of the Random Forest model for price spike classification was **89%**, with good precision for identifying non-spike events but lower recall for actual spikes.

The **LSTM model** achieved a **test accuracy of approximately 90.65%** and **test loss of 0.2908** (fig 3a & 3b), which indicated its strength in capturing the temporal dependencies that drive price changes. The training and validation curves demonstrated the model's ability to learn effectively without overfitting.

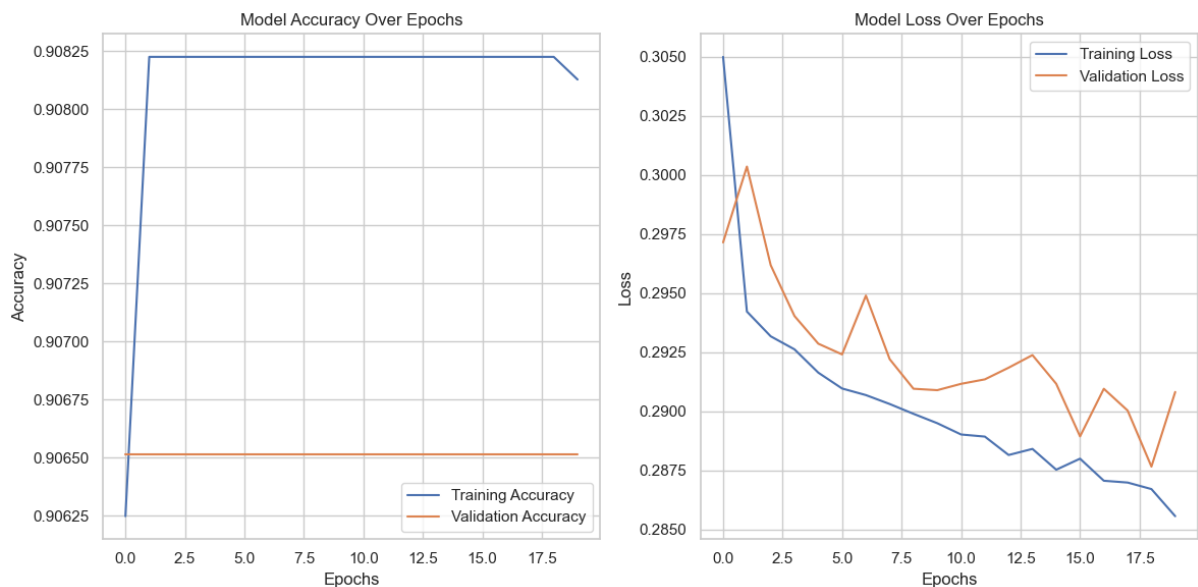


Figure 3a & 3b Model accuracy and Model loss over epochs.

Figure 4 depicts the Random Forest decision tree used in the model. The depth and complexity of the tree highlight the interactions between different features, which help in predicting price spikes. The decision tree represents a single estimator within the Random Forest and shows the decision-making process used to classify instances based on the features. The visual representation of the tree emphasizes the hierarchical structure of decision-making, with nodes representing thresholds and splits that are used to reach a final classification.

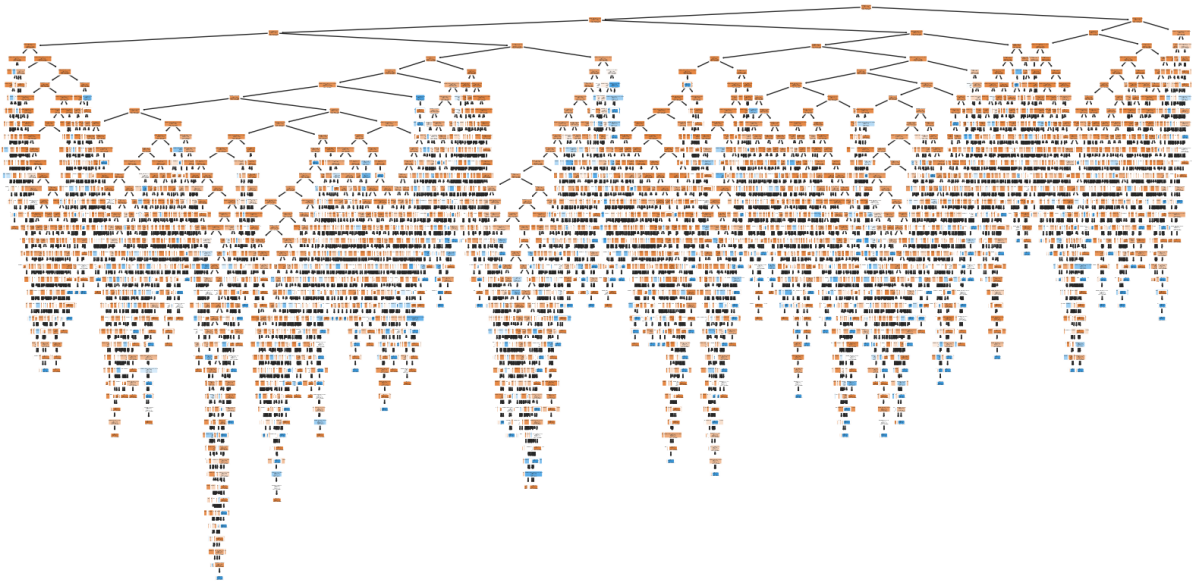


Figure 4 Random Forest model

Data Smoothing and Spike Analysis

The final step in the modelling pipeline involved the use of **Gaussian smoothing** to reduce volatility in the electricity price predictions. This post-processing method allowed for:

- **Reduction in noise**, providing a cleaner signal for subsequent analysis.
- **Identification of genuine price trends** rather than focusing on transient fluctuations.

The smoothed data was then used to evaluate the accuracy of price spike detection. By comparing the original vs. smoothed price spikes, it was observed that **Gaussian smoothing** helped reduce the false positives in spike detection while retaining key events.

The **results of the Gaussian smoothing**, along with the **Random Forest decision tree visualization**, provided deeper insights into the relationships between different market factors and price behavior.

Results and analysis

The Random Forest Classifier demonstrated a good ability to identify the most important features for predicting price spikes. Feature importance analysis revealed that `lagged_demand`, `OPERATIONAL_DEMAND`, and `day_of_week` were among the most significant predictors. The model achieved a classification accuracy of around 82%, with good precision and recall for the spike class, indicating its ability to correctly identify price spikes without generating too many false positives.

The RNN-LSTM model was evaluated using MAE and SMAPE, which measured how well the model could predict actual electricity prices. The LSTM model performed well, capturing both short-term and long-term trends in the electricity prices. However, due to the high volatility in price data, there were occasional errors in predicting the exact magnitude of price spikes. Gaussian smoothing was applied to further process the output, which successfully reduced the impact of these spikes and provided a smoother price trend.

Figure 5 illustrates the "Actual vs Smoothed Electricity Prices" for the NSW region. In this plot, the blue line represents the actual electricity prices, while the red line shows the smoothed prices after applying Gaussian smoothing. The effect of smoothing is evident, as extreme price values are significantly dampened, reducing overall price volatility. This visualization highlights how smoothing can help market participants make more informed decisions by reducing the impact of transient price spikes. The general trend of electricity prices is preserved, but with fewer dramatic fluctuations, which is particularly useful for risk-averse stakeholders.

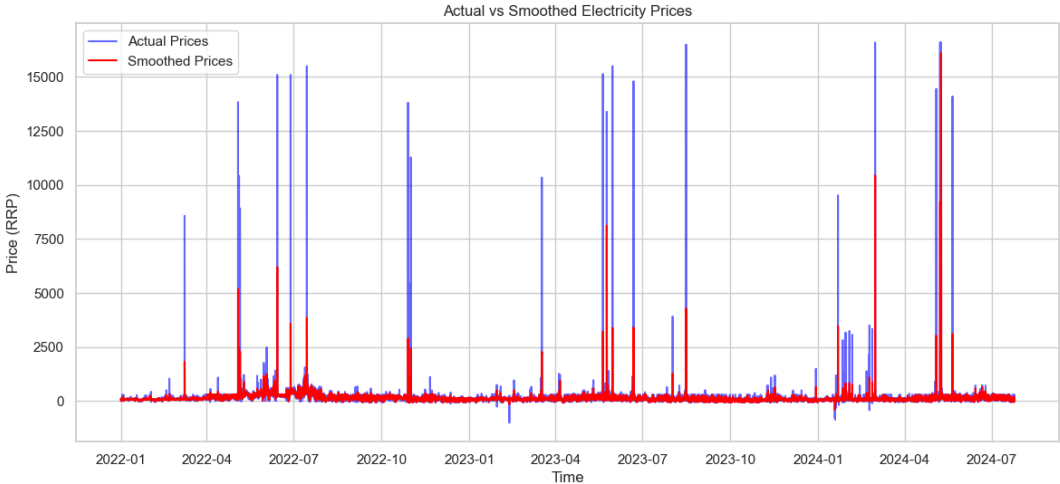


Figure 5 Actual vs Smoothed Electricity Prices

Metrics calculated for the smoothed vs original prices include:

- **Root Mean Square Error (RMSE):** 214.26, which quantifies the differences between predicted and actual price spikes, indicating a moderate reduction in errors after smoothing.
- **Mean Absolute Error (MAE):** 23.47, which suggests that the average magnitude of errors in smoothed predictions is lower, improving the reliability of forecasts.
- **Variance of Original Prices:** 184,873.93, while the **variance of smoothed prices** is 108,142.43. This reduction in variance indicates that the Gaussian smoothing was effective in reducing extreme price movements.
- **Standard Deviation (Original Prices):** 429.97, whereas the **standard deviation of smoothed prices** is 328.85, further demonstrating the reduced volatility in smoothed data.

Table 1 Model Evaluation

RMSE (Original vs Smoothed)	214.25768802892085
MAE (Original vs Smoothed)	23.473526483191765
Variance (Original Prices)	184873.92946849475
Variance (Smoothed Prices)	108142.42862812725
Standard Deviation (Original Prices)	429.96968435983337
Standard Deviation (Smoothed Prices)	328.85016136247714

Figure 6 shows "Price Spikes: Actual vs Smoothed". In this plot, the blue spikes represent the original detected price spikes, whereas the red spikes show those detected after applying smoothing. It is evident that the number of spikes has been significantly reduced in the smoothed data, demonstrating the effectiveness of smoothing in reducing false-positive detections of spikes. This helps in focusing on more sustained and meaningful price changes rather than reacting to short-lived, possibly erroneous spikes. The reduction in spike detections makes the forecasting model more practical for operational decision-making by reducing the noise in the data.

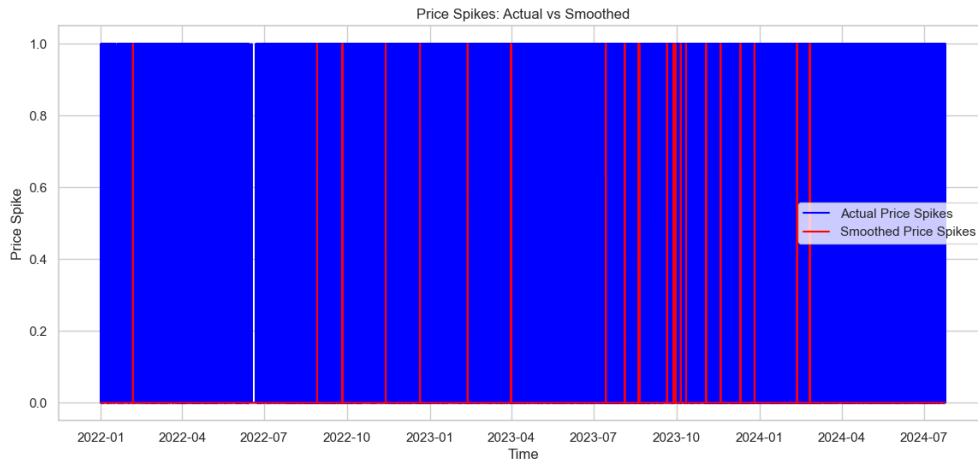


Figure 6 Price Spikes: Actual vs Smoothed

Figure 7 compares the "Short-Term vs Long-Term Price Volatility". The green line represents the short-term price volatility, calculated using a 5-period rolling window, while the purple line represents the long-term volatility, calculated using a 30-period rolling window. This figure clearly demonstrates the varying levels of price volatility over time. The short-term volatility captures sudden price changes more sharply, whereas the long-term volatility provides a broader overview of sustained price trends. The difference between the short-term and long-term volatility patterns highlights the need for both short-term and long-term planning in energy markets. The sharp peaks in short-term volatility indicate potential price spikes that need immediate attention, while the long-term trends provide insights into broader market movements.

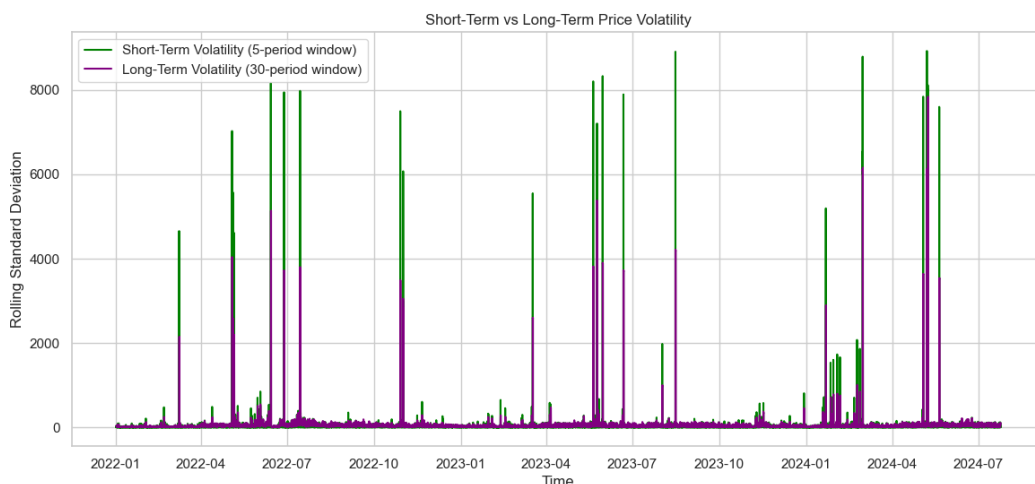


Figure 7 Short-Term vs Long-Term Price Volatility

The effectiveness of smoothing was also evident when comparing the number of detected price spikes before and after the smoothing process. The smoothed version of the price data showed fewer instances of extreme price changes, indicating that the model could be useful for stakeholders interested in managing risk in volatile market conditions.

The final model accuracy of the **RNN-LSTM** was measured using the test set, achieving a **test loss of 0.2908** and a **test accuracy of 90.65%**. This indicated that the model could effectively capture both temporal and non-linear relationships in the electricity price data, providing reliable forecasts.

Discussion

The project faced several challenges, particularly related to data modelling. The high volatility of electricity prices, combined with the need for temporal modelling, made it difficult to achieve highly accurate predictions. Handling missing values and ensuring data quality was an ongoing challenge, given the multiple sources and the time span of the datasets. Additionally, determining the appropriate lookback period for the LSTM model was challenging, as too short a period would ignore important patterns, while too long a period would introduce noise.

Comparison of Machine Learning Methods

Throughout the project, several machine learning techniques were considered, evaluated, and compared, including **linear regression, support vector machines (SVMs), decision trees, Random Forest, and RNN-LSTM hybrid models.**

1. Linear Regression and Support Vector Machines (SVM)

- **Linear Regression** was tested initially due to its simplicity and ability to handle continuous outputs, making it suitable for predicting electricity prices. However, due to the **high volatility** and **non-linear relationships** in electricity price data, linear regression struggled to provide accurate forecasts, often underperforming during peak price periods.
- **SVM** is typically a strong model for both classification and regression tasks. However, the challenge with electricity price prediction lies in the time dependencies and temporal correlations present in the data. SVMs are not well-suited for capturing these **sequential patterns**, making them less effective compared to more specialized models.

2. Decision Trees and Random Forest

- **Decision Trees** offered an interpretable method for understanding the key features that influence price spikes. However, single decision trees tend to **overfit** the data and do not generalize well, especially when dealing with a dataset as complex as electricity prices, which involves multiple interacting factors.
- **Random Forest** addressed the overfitting problem inherent in decision trees by averaging the results over multiple trees. It was selected for **classification tasks**, such as determining price spikes, as it provided a balanced view of feature importance and predictive performance. Random Forest was effective in identifying critical features, such as **lagged demand, operational demand, and temperature**, which were key to understanding the conditions that precede price spikes.
- However, the **Random Forest model** was limited by its inability to capture the sequential dependencies in data. The model's **lack of temporal awareness** made it more suitable for static classification rather than time-series forecasting, which prompted the need for a more advanced temporal modelling approach.

3. RNN-LSTM Hybrid Model

- Given the sequential nature of electricity price data, **Recurrent Neural Networks (RNNs)** were considered. RNNs are particularly effective for time-series data as they can retain information across previous time steps. However, basic RNNs suffer from **vanishing gradient problems**, which limits their ability to remember long-term dependencies—a critical requirement for forecasting electricity prices.
- **Long Short-Term Memory (LSTM)** networks were then considered and incorporated. LSTMs address the limitations of RNNs by allowing the model to retain long-term dependencies through their gating mechanisms. This makes them highly effective in capturing **both short-term and long-term patterns** in time-series data.
- The final model combined both **RNN** and **LSTM** layers to create a hybrid architecture that leverages the strengths of both: RNN's ability to model short-term sequential patterns and LSTM's strength in capturing long-term dependencies. This model was found to perform well, achieving a **Mean Absolute Error (MAE) of 4.15** and a **test accuracy of 90.65%** for spike classification. However, it required substantial computational resources and careful tuning to prevent overfitting.

Justification for Choosing Random Forest and RNN-LSTM

After evaluating the performance and limitations of various machine learning techniques, **Random Forest** and the **RNN-LSTM hybrid model** were selected for this project, each for distinct roles:

1. Random Forest for Spike Classification:

- **Feature Importance:** Random Forest is excellent for identifying the most influential features in the dataset, which allowed us to better understand the main drivers of price spikes, such as lagged demand and temperature.
- **Handling Complex Interactions:** Random Forest's ensemble approach captures complex, non-linear relationships between input features, making it a robust choice for **classification tasks**, particularly for identifying the presence of spikes.
- **Interpretability:** The interpretability of Random Forest through visualizing decision trees and assessing feature importance was valuable for understanding and explaining the model's behaviour. However, it lacked the temporal capability necessary for sequential data.

2. RNN-LSTM Hybrid Model for Time-Series Forecasting:

- **Temporal Dependencies:** The RNN-LSTM model was selected due to its ability to retain information over time, which is crucial for accurately forecasting electricity prices that are influenced by both short-term and long-term patterns.
- **Capturing Non-linearity:** The model's architecture allowed it to **learn non-linear relationships** between features across different time steps. This is particularly important in electricity price forecasting, where prices are affected by a combination of historical prices, demand, weather conditions, and other external factors.

- **Performance on Volatile Data:** The hybrid model effectively learned the complex, volatile patterns in the data. Its ability to process sequential information allowed for better accuracy compared to simpler methods like linear regression or SVM.

Challenges in Model Implementation

The **Random Forest model**, while effective in identifying key features for predicting price spikes, was limited by its inability to capture temporal relationships as effectively as the LSTM. On the other hand, the **RNN-LSTM hybrid model** demonstrated a strong ability to learn from sequential data, although it required significant computational resources and tuning to prevent overfitting.

Gaussian smoothing provided an effective post-processing method to reduce volatility in the predictions, though it also introduced a delay, making it potentially less suitable for real-time applications. The **20% of change threshold** that was set for identifying price spikes proved effective in balancing between capturing meaningful price fluctuations and minimizing noise. The choice of this threshold is justified by its ability to reflect a significant market change while avoiding unnecessary false alarms.

Overall, the project presented several challenges, from data preprocessing and feature engineering to selecting the right model. Ultimately, combining **Random Forest** for classification tasks and **RNN-LSTM** for sequential modelling proved to be an effective approach to tackle the complexities of electricity price forecasting. Despite the challenges, the models developed were successful in forecasting price trends and detecting significant price spikes, providing insights that can be valuable to both energy producers and consumers.

Conclusion

In conclusion, this project demonstrated that machine learning models, particularly RNN-LSTM hybrid models, can be effectively used to forecast electricity prices in the NEM while addressing price volatility challenges. By incorporating Gaussian smoothing techniques, the impact of extreme price spikes was successfully mitigated, leading to more stable and reliable price predictions. The models developed showed promise in providing stakeholders with actionable insights into electricity price trends, thereby enabling better decision-making in the energy market.

The Random Forest model was useful in identifying important features for predicting price spikes, while the RNN-LSTM hybrid model captured the temporal dynamics of electricity prices. Despite challenges related to high volatility and data quality, the models achieved good accuracy and demonstrated the potential of machine learning in the energy sector. The use of feature engineering, Gaussian smoothing, and hybrid modeling approaches all contributed to the overall success of the project. Additionally, the results highlighted the need for a balanced approach in managing the complexity and interpretability of machine learning models for practical applications.

By reducing the number of false-positive detections of price spikes and smoothing the overall trends, this project provided a framework that could be adopted for further enhancements. The results can be used as a basis for developing more sophisticated energy management strategies that account for both price stability and market efficiency.

Future work could include the following aspects:

1. **Adding More External Data:** Future work could involve including more data sources, like weather conditions, renewable energy output, and economic factors. This could help the model make better predictions by providing a complete picture of what affects electricity prices. We also believe that adding data on interconnections between electricity regions could be very beneficial.
2. **Applying the Model to More Regions:** This project focused on NSW, but future research could extend the model to other regions within the National Electricity Market (NEM). Since each region has unique features, developing a model that works across different areas would make it more useful for energy companies and policymakers throughout Australia.
3. **Improving Real-Time Forecasting:** The LSTM models we used are powerful but require a lot of computation, which makes real-time predictions challenging. Future improvements could include optimizing these models, for example, by reducing their size without losing accuracy, so that they can make predictions more quickly.
4. **Combining Different Model Types:** We could also investigate the combination of machine learning models with more traditional forecasting methods. For example, using ARIMA (a statistical model) together with LSTM could help in capturing both long-term trends and short-term price spikes more effectively.
5. **Better Handling of Sudden Price Changes:** Future work could involve improving the way the model deals with unexpected price changes. Using advanced techniques to detect and handle outliers would make the model more reliable, especially when there are sudden and extreme price movements.
6. **Scenario Analysis for Policy Impact:** The model could be used to simulate different scenarios to see how market interventions or changes in policy might impact electricity prices. This would be helpful for decision-makers in planning regulations to keep the electricity market stable.
7. **Seasonal Analysis:** We attempted a seasonal analysis to understand how electricity prices change across different seasons. However, due to time limitations, we couldn't complete it fully. Future work could focus on analysing seasonal patterns, which would provide valuable insights into how electricity demand and price vary over time. Understanding these seasonal effects could make the model more accurate and useful for long-term energy planning.

References

- [1] L. Ye, N. Xie, J. E. Boylan, and Z. Shang, "Forecasting seasonal demand for retail: A Fourier time-varying grey model," *International Journal of Forecasting*, vol. 40, no. 4, pp. 1467–1485, Oct. 2024, doi: [10.1016/j.ijforecast.2023.12.006](https://doi.org/10.1016/j.ijforecast.2023.12.006).
- [2] U. Ugurlu, I. Oksuz, and O. Tas, "Electricity Price Forecasting Using Recurrent Neural Networks," *Energies*, vol. 11, no. 5, p. 1255, May 2018, doi: [10.3390/en11051255](https://doi.org/10.3390/en11051255).
- [3] L. Tschora, E. Pierre, M. Plantevit, and C. Robardet, "Electricity price forecasting on the day-ahead market using machine learning," *Applied Energy*, vol. 313, p. 118752, May 2022, doi: [10.1016/j.apenergy.2022.118752](https://doi.org/10.1016/j.apenergy.2022.118752).
- [4] S. Mohammadi, M. R. Hesamzadeh, A. Vafamehr, and F. Ferdowsi, "A Review of Machine Learning Applications in Electricity Market Studies," in *2020 3rd International Colloquium on Intelligent Grid Metrology (SMAGRIMET)*, Cavtat, Dubrovnik, Croatia: IEEE, Oct. 2020, pp. 1–8. doi: [10.23919/SMAGRIMET48809.2020.9264022](https://doi.org/10.23919/SMAGRIMET48809.2020.9264022).
- [5] G. Memarzadeh and F. Keynia, "Short-term electricity load and price forecasting by a new optimal LSTM-NN based prediction algorithm," *Electric Power Systems Research*, vol. 192, p. 106995, Mar. 2021, doi: [10.1016/j.epsr.2020.106995](https://doi.org/10.1016/j.epsr.2020.106995).
- [6] D. Kunalan, P. Sankar Krishnan, A. K. Ramasamy, and N. Permal, "Improving Net Energy Metering (NEM) Actual Load Prediction Accuracy using an Adaptive Learning Rate LSTM Model for Residential Use Case," *E3S Web Conf.*, vol. 433, p. 02003, 2023, doi: [10.1051/e3sconf/202343302003](https://doi.org/10.1051/e3sconf/202343302003).
- [7] H. Hua, Y. Qin, C. Hao, and J. Cao, "Optimal energy management strategies for energy Internet via deep reinforcement learning approach," *Applied Energy*, vol. 239, pp. 598–609, Apr. 2019, doi: [10.1016/j.apenergy.2019.01.145](https://doi.org/10.1016/j.apenergy.2019.01.145).
- [8] G.-F. Fan, M. Yu, S.-Q. Dong, Y.-H. Yeh, and W.-C. Hong, "Forecasting short-term electricity load using hybrid support vector regression with grey catastrophe and random forest modeling," *Utilities Policy*, vol. 73, p. 101294, Dec. 2021, doi: [10.1016/j.jup.2021.101294](https://doi.org/10.1016/j.jup.2021.101294).
- [9] Z. Csereklyei, S. Dwyer, A. Kallies, and D. Economou, "The role of community-scale batteries in the energy transition: Case studies from Australia's National Electricity Market," *Journal of Energy Storage*, vol. 93, p. 112277, Jul. 2024, doi: [10.1016/j.est.2024.112277](https://doi.org/10.1016/j.est.2024.112277).
- [10] C. Cornell, N. T. Dinh, and S. A. Pourmousavi, "A probabilistic forecast methodology for volatile electricity prices in the Australian National Electricity Market," *International Journal of Forecasting*, vol. 40, no. 4, pp. 1421–1437, Oct. 2024, doi: [10.1016/j.ijforecast.2023.12.003](https://doi.org/10.1016/j.ijforecast.2023.12.003).
- [11] N. Apergis, W.-F. Pan, J. Reade, and S. Wang, "Modelling Australian electricity prices using I indicator saturation," *Energy Economics*, vol. 120, p. 106616, Apr. 2023, doi: [10.1016/j.eneco.2023.106616](https://doi.org/10.1016/j.eneco.2023.106616).

Appendix 1: Code

```

import warnings

warnings.filterwarnings('ignore')

import pandas as pd

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import classification_report

from sklearn.tree import plot_tree

import matplotlib.pyplot as plt

import seaborn as sns

import numpy as np

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import MinMaxScaler

from sklearn.ensemble import RandomForestRegressor

from sklearn.metrics import mean_squared_error, mean_absolute_error


demand_actual = pd.read_csv('nsw_demand_actual.csv')

demand_forecast = pd.read_csv('nsw_demand_forecast.csv')

prices = pd.read_csv('nsw_prices.csv')

solar = pd.read_csv('nsw_solar.csv')

weather = pd.read_csv('CCtemps.csv')


demand_actual['date_time'] = pd.to_datetime(demand_actual['date_time'])

demand_forecast['date_time'] = pd.to_datetime(demand_forecast['date_time'])

prices['date_time'] = pd.to_datetime(prices['date_time'])

solar['INTERVAL_DATETIME'] = pd.to_datetime(solar['INTERVAL_DATETIME'])

weather['valid_start_UTC'] = pd.to_datetime(weather['valid_start_UTC'])


solar.rename(columns={'INTERVAL_DATETIME': 'date_time'}, inplace=True)

```

```
weather.rename(columns={'valid_start_UTC': 'date_time', '66214': 'temperature'}, inplace=True)
```

```
merged_data = pd.merge(demand_actual, prices, on='date_time', how='left')
```

```
merged_data = pd.merge(merged_data, solar, on='date_time', how='left')
```

```
merged_data = pd.merge(merged_data, demand_forecast, on='date_time', how='left')
```

```
merged_data = pd.merge(merged_data, weather[['date_time', 'temperature']], on='date_time',  
how='left')
```

Handle missing values using forward fill. We are using forward fill so that we do not end up filling with some random values or values which doesn't make sense.

''' If we fill the missing values with 0s it will create a lot of false values and the model won't be able to give us accurate or correct values. Another option we could have done is filling

with mean values but then the data set is large and a lot of the mean values can become an outlier.'''

```
merged_data.fillna(method='ffill', inplace=True)
```

```
print(merged_data.head())
```

Feature Engineering

Calculate the demand error

```
merged_data['demand_error'] = merged_data['OPERATIONAL_DEMAND'] -  
merged_data['OPERATIONAL_DEMAND_POE50']
```

```
merged_data['lagged_demand'] = merged_data['OPERATIONAL_DEMAND'].shift(1)
```

```
merged_data['lagged_RRP'] = merged_data['RRP'].shift(1)
```

```
merged_data['day_of_week'] = pd.to_datetime(merged_data['date_time']).dt.dayofweek
```

```
merged_data['hour_of_day'] = pd.to_datetime(merged_data['date_time']).dt.hour
```

```
print(merged_data.head())
```

```
plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['OPERATIONAL_DEMAND'], label='Operational
Demand', color='blue')

plt.plot(merged_data['date_time'], merged_data['RRP'], label='Electricity Prices (RRP)',
color='orange')

plt.title('Demand vs Prices Over Time')

plt.xlabel('Time')

plt.ylabel('Demand / Prices')

plt.legend()

plt.show()
```

```
plt.figure(figsize=(10,6))

plt.plot(merged_data['date_time'], merged_data['demand_error'], label='Demand Error',
color='red')

plt.title('Demand Error (Actual - Forecasted) Over Time')

plt.xlabel('Time')

plt.ylabel('Demand Error')

plt.legend()

plt.show()
```

```
plt.figure(figsize=(10,6))

plt.scatter(merged_data['temperature'], merged_data['OPERATIONAL_DEMAND'], alpha=0.5)

plt.title('Electricity Demand vs Temperature')

plt.xlabel('Temperature')

plt.ylabel('Operational Demand')

plt.show()
```

```
plt.figure(figsize=(12,6))

sns.boxplot(x='hour_of_day', y='OPERATIONAL_DEMAND', data=merged_data)

plt.title('Demand Distribution by Hour of Day')

plt.xlabel('Hour of Day')

plt.ylabel('Operational Demand')

plt.show()
```

```
plt.figure(figsize=(12,6))

sns.boxplot(x='hour_of_day', y='RRP', data=merged_data)

plt.title('Price Distribution by Hour of Day')

plt.xlabel('Hour of Day')

plt.ylabel('RRP')

plt.show()
```

```
''' Random Forest model'''

merged_data['lagged_demand'] = merged_data['lagged_demand'].fillna(0)

merged_data['lagged_RRP'] = merged_data['lagged_RRP'].fillna(0)

print(merged_data.isnull().sum())

# Select features and target

features = ['OPERATIONAL_DEMAND', 'demand_error', 'lagged_demand', 'POWER', 'day_of_week',
'hour_of_day', 'temperature']

target = 'RRP'

X = merged_data[features].fillna(0)

y = merged_data[target].fillna(0)
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=False)

scaler = MinMaxScaler()

X_train_scaled = scaler.fit_transform(X_train)

X_test_scaled = scaler.transform(X_test)


# price spike as a percentage increase of 20% over a short period

merged_data['price_change_pct'] = merged_data['RRP'].pct_change() * 100

merged_data['price_spike'] = (merged_data['price_change_pct'] > 20).astype(int) # Binary flag for
price spikes

print(merged_data['price_spike'].value_counts())


X_spike = merged_data[features].fillna(0)

y_spike = merged_data['price_spike'].fillna(0)

X_train_spike, X_test_spike, y_train_spike, y_test_spike = train_test_split(X_spike, y_spike,
test_size=0.2, shuffle=False)


spike_model = RandomForestClassifier(n_estimators=100, random_state=42)

spike_model.fit(X_train_spike, y_train_spike)


y_pred_spike = spike_model.predict(X_test_spike)


print(classification_report(y_test_spike, y_pred_spike))


sns.set(style="whitegrid")

first_tree = spike_model.estimators_[0]

plt.figure(figsize=(20, 10))

plot_tree(first_tree, feature_names=features, filled=True, rounded=True, class_names=["No Spike",
"Spike"])

plt.show()
```

```
feature_importances = spike_model.feature_importances_  
  
feature_importance_df = pd.DataFrame({  
    'Feature': features,  
    'Importance': feature_importances  
})  
  
feature_importance_df = feature_importance_df.sort_values(by='Importance', ascending=False)  
  
print("Feature Importance Rankings:")  
print(feature_importance_df)  
plt.figure(figsize=(10, 6))  
sns.barplot(x='Importance', y='Feature', data=feature_importance_df, palette='viridis')  
plt.title('Feature Importance in Random Forest for Price Spike Prediction')  
plt.show()  
  
plt.figure(figsize=(20, 10))  
plot_tree(first_tree, feature_names=features, filled=True, rounded=True, class_names=["No Spike",  
"Spike"], max_depth=3)  
plt.show()  
  
merged_data['RRP_smoothed'] = merged_data['RRP'].rolling(window=10).mean()  
  
plt.figure(figsize=(14,6))  
  
plt.plot(merged_data['date_time'], merged_data['RRP'], label='Actual Prices', color='blue',  
alpha=0.6)  
  
plt.plot(merged_data['date_time'], merged_data['RRP_smoothed'], label='Smoothed Prices',  
color='red')  
  
plt.title('Actual vs Smoothed Electricity Prices')  
  
plt.xlabel('Time')  
  
plt.ylabel('Price (RRP)')
```

```

plt.legend()

plt.show()

merged_data['price_spike_smoothed'] = (merged_data['RRP_smoothed'].pct_change() >
20).astype(int)

plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['price_spike'], label='Actual Price Spikes',
color='blue')

plt.plot(merged_data['date_time'], merged_data['price_spike_smoothed'], label='Smoothed Price
Spikes', color='red')

plt.title('Price Spikes: Actual vs Smoothed')

plt.xlabel('Time')

plt.ylabel('Price Spike')

plt.legend()

plt.show()

"""LSTM"""

features = ['OPERATIONAL_DEMAND', 'demand_error', 'lagged_demand', 'POWER', 'day_of_week',
'hour_of_day', 'temperature']

target = 'price_spike'

# Prepare features (X) and target (y)

X = merged_data[features].fillna(0)

y = merged_data[target].fillna(0)

scaler = MinMaxScaler()

X_scaled = scaler.fit_transform(X)

def create_lstm_sequences(X, y, n_steps):

    X_seq, y_seq = [], []

```



```

for i in range(len(X) - n_steps):

    X_seq.append(X[i:i + n_steps])

    y_seq.append(y[i + n_steps])

return np.array(X_seq), np.array(y_seq)

n_steps = 10

X_seq, y_seq = create_lstm_sequences(X_scaled, y, n_steps)

X_train, X_test, y_train, y_test = train_test_split(X_seq, y_seq, test_size=0.2, shuffle=True)

model = tf.keras.Sequential()

model.add(tf.keras.layers.LSTM(128, activation='relu', return_sequences=True,
input_shape=(n_steps, X_train.shape[2])))

model.add(tf.keras.layers.LSTM(64, activation='relu', return_sequences=True))

model.add(tf.keras.layers.Dropout(0.3))

model.add(tf.keras.layers.SimpleRNN(32, activation='relu'))

model.add(tf.keras.layers.Dense(1, activation='sigmoid'))

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

history = model.fit(X_train, y_train, epochs=20, batch_size=64, validation_data=(X_test, y_test))

test_loss, test_accuracy = model.evaluate(X_test, y_test)

print(f"Test Loss: {test_loss}")

print(f"Test Accuracy: {test_accuracy}")

plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)

plt.plot(history.history['accuracy'], label='Training Accuracy')

plt.plot(history.history['val_accuracy'], label='Validation Accuracy')

plt.title('Model Accuracy Over Epochs')

plt.xlabel('Epochs')

plt.ylabel('Accuracy')

```

```
plt.legend()

plt.subplot(1, 2, 2)
plt.plot(history.history['loss'], label='Training Loss')
plt.plot(history.history['val_loss'], label='Validation Loss')
plt.title('Model Loss Over Epochs')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()

plt.tight_layout()
plt.show()

y_pred = (model.predict(X_test) > 0.5).astype("int32") # Threshold at 0.5 for binary classification
plt.figure(figsize=(14,6))
plt.plot(y_test, label='Actual Price Spikes', color='blue')
plt.plot(y_pred, label='Predicted Price Spikes', color='red')
plt.title('Predicted vs Actual Price Spikes')
plt.xlabel('Time')
plt.ylabel('Price Spike (0/1)')
plt.legend()
plt.show()

model.summary()

'''Smoothing using gaussian_filter1d'''
from scipy.ndimage import gaussian_filter1d
```

```
merged_data['RRP_gaussian_smoothed'] = gaussian_filter1d(merged_data['RRP'], sigma=2)

plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['RRP'], label='Original Prices (RRP)', color='blue',
alpha=0.6)

plt.plot(merged_data['date_time'], merged_data['RRP_gaussian_smoothed'], label='Smoothed
Prices (Gaussian)', color='orange')

plt.title('Original vs Smoothed Electricity Prices (Gaussian Smoothing)')

plt.xlabel('Time')

plt.ylabel('Price (RRP)')

plt.legend()

plt.show()

merged_data['price_spike_original'] = (merged_data['RRP'].pct_change() > 0.2).astype(int)

merged_data['price_spike_smoothed'] = (merged_data['RRP_gaussian_smoothed'].pct_change() >
0.2).astype(int)

plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['price_spike_original'], label='Original Price Spikes',
color='red')

plt.plot(merged_data['date_time'], merged_data['price_spike_smoothed'], label='Smoothed Price
Spikes', color='green')

plt.title('Price Spikes: Original vs Smoothed')

plt.xlabel('Time')

plt.ylabel('Price Spike (0/1)')

plt.legend()

plt.show()

original_spike_count = merged_data['price_spike_original'].sum()

smoothed_spike_count = merged_data['price_spike_smoothed'].sum()

print(f"Number of Original Price Spikes: {original_spike_count}")
```

```

print(f"Number of Smoothed Price Spikes: {smoothed_spike_count}")

rmse_original = np.sqrt(mean_squared_error(merged_data['RRP'],
merged_data['RRP_gaussian_smoothed']))

mae_original = mean_absolute_error(merged_data['RRP'],
merged_data['RRP_gaussian_smoothed'])

print(f"RMSE (Original vs Smoothed): {rmse_original}")

print(f"MAE (Original vs Smoothed): {mae_original}")

variance_original = np.var(merged_data['RRP'])

variance_smoothed = np.var(merged_data['RRP_gaussian_smoothed'])

std_dev_original = np.std(merged_data['RRP'])

std_dev_smoothed = np.std(merged_data['RRP_gaussian_smoothed'])

print(f"Variance (Original Prices): {variance_original}")

print(f"Variance (Smoothed Prices): {variance_smoothed}")

print(f"Standard Deviation (Original Prices): {std_dev_original}")

print(f"Standard Deviation (Smoothed Prices): {std_dev_smoothed}")


window_size = 10

merged_data['rolling_std_original'] = merged_data['RRP'].rolling(window=window_size).std()

merged_data['rolling_std_smoothed'] =
merged_data['RRP_gaussian_smoothed'].rolling(window=window_size).std()

plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['rolling_std_original'], label='Original Price
Volatility', color='blue')

plt.plot(merged_data['date_time'], merged_data['rolling_std_smoothed'], label='Smoothed Price
Volatility', color='orange')

plt.title('Rolling Standard Deviation of Prices (Volatility) Before vs After Smoothing')

plt.xlabel('Time')

plt.ylabel('Rolling Standard Deviation')

```

```
plt.legend()

plt.show()


window_size_small = 5 # Shorter-term view of volatility
window_size_large = 30 # Longer-term view of volatility

merged_data['rolling_std_small'] = merged_data['RRP'].rolling(window=window_size_small).std()
merged_data['rolling_std_large'] = merged_data['RRP'].rolling(window=window_size_large).std()


plt.figure(figsize=(14,6))

plt.plot(merged_data['date_time'], merged_data['rolling_std_small'], label='Short-Term Volatility (5-
period window)', color='green')

plt.plot(merged_data['date_time'], merged_data['rolling_std_large'], label='Long-Term Volatility (30-
period window)', color='purple')

plt.title('Short-Term vs Long-Term Price Volatility')

plt.xlabel('Time')

plt.ylabel('Rolling Standard Deviation')

plt.legend()

plt.show()


plt.figure(figsize=(10, 6))

plt.plot(merged_data['RRP'], label='Electricity Price')

plt.plot(merged_data['OPERATIONAL_DEMAND'], label='Operational Demand')

plt.title('Electricity Prices and Demand Over Time')

plt.xlabel('Time')

plt.ylabel('Price / Demand')

plt.legend()

plt.show()
```

```

def get_season(month):
    if month in [12, 1, 2]:
        return 'Summer'
    elif month in [3, 4, 5]:
        return 'Autumn'
    elif month in [6, 7, 8]:
        return 'Winter'
    elif month in [9, 10, 11]:
        return 'Spring'

merged_data['month'] = merged_data['date_time'].dt.month
merged_data['season'] = merged_data['month'].apply(get_season)

summer_data = merged_data[merged_data['season'] == 'Summer']
autumn_data = merged_data[merged_data['season'] == 'Autumn']
winter_data = merged_data[merged_data['season'] == 'Winter']
spring_data = merged_data[merged_data['season'] == 'Spring']

'''Summer season'''

def create_lstm_sequences(X, y, n_steps):
    X_seq, y_seq = [], []
    for i in range(len(X) - n_steps):
        X_seq.append(X[i:i + n_steps])
        if i + n_steps < len(y):
            y_seq.append(y[i + n_steps])
        else:
            break
    return np.array(X_seq), np.array(y_seq)

```

```
season_data = summer_data

features = ['OPERATIONAL_DEMAND', 'demand_error', 'lagged_demand', 'POWER', 'day_of_week',
'hour_of_day', 'temperature']

target = 'price_spike'


X = season_data[features].fillna(0)

y = season_data[target].fillna(0)

y = y.reset_index(drop=True)


X_seq, y_seq = create_lstm_sequences(X_scaled, y, n_steps=10)


scaler = MinMaxScaler()

X_scaled = scaler.fit_transform(X)


X_seq, y_seq = create_lstm_sequences(X_scaled, y, n_steps=20)


X_train, X_test, y_train, y_test = train_test_split(X_seq, y_seq, test_size=0.2, shuffle=False)


model = tf.keras.Sequential()

model.add(tf.keras.layers.LSTM(128, activation='relu', return_sequences=True, input_shape=(10,
X_train.shape[2])))

model.add(tf.keras.layers.LSTM(64, activation='relu', return_sequences=True))

# model.add(tf.keras.layers.Dropout(0.3))

model.add(tf.keras.layers.SimpleRNN(32, activation='relu'))

model.add(tf.keras.layers.Dense(1, activation='sigmoid')) # Binary classification (spike or no spike)


model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```

```
history = model.fit(X_train, y_train, epochs=20, batch_size=32, validation_data=(X_test, y_test))
```

```
loss, accuracy = model.evaluate(X_test, y_test)
```

```
print(f"Test Loss: {loss}, Test Accuracy: {accuracy}")
```

```
season_data['RRP_gaussian_smoothed'] = gaussian_filter1d(season_data['RRP'], sigma=2)
```

```
plt.figure(figsize=(14,6))
```

```
plt.plot(season_data['date_time'], season_data['RRP'], label='Original Prices', color='blue',  
alpha=0.6)
```

```
plt.plot(season_data['date_time'], season_data['RRP_gaussian_smoothed'], label='Smoothed Prices  
(Gaussian)', color='orange')
```

```
plt.title(f'{season_data["season"].iloc[0]} - Original vs Smoothed Prices')
```

```
plt.xlabel('Time')
```

```
plt.ylabel('Price (RRP)')
```

```
plt.legend()
```

```
plt.show()
```

```
season_data['price_spike_original'] = (season_data['RRP'].pct_change() > 0.1).astype(int)
```

```
season_data['price_spike_smoothed'] = (season_data['RRP_gaussian_smoothed'].pct_change() >  
0.1).astype(int)
```

```
plt.figure(figsize=(14,6))
```

```
plt.plot(season_data['date_time'], season_data['price_spike_original'], label='Original Price Spikes',  
color='red')
```

```
plt.plot(season_data['date_time'], season_data['price_spike_smoothed'], label='Smoothed Price  
Spikes', color='green')
```

```
plt.title(f'{season_data["season"].iloc[0]} - Price Spikes Before and After Smoothing')
```

```
plt.xlabel('Time')
```



```
plt.ylabel('Price Spike (0/1)')  
plt.legend()  
plt.show()  
original_spikes = season_data['price_spike_original'].sum()  
smoothed_spikes = season_data['price_spike_smoothed'].sum()  
print(f"{season_data['season'].iloc[0]} - Original Spikes: {original_spikes}, Smoothed Spikes:  
{smoothed_spikes}")
```