

Utilisation of Time Series Data for Forecasting Oil Temperature and Countries Exchange Rates

BENNABI Malek Mohamed
Université Paris Cité

April 24, 2024

Abstract

In this paper, we harness the power of time series data to forecast two critical parameters: oil temperature and exchange rates. Our research is grounded in the analysis of electricity consumption patterns, which serve as a rich source of time series data. We employ machine learning techniques to identify patterns and trends in this data, which are then used to predict oil temperature and exchange rates.

Our approach involves forecasting the Oil temperature within the Electricity Transformer Temperature data. We also conduct a comparative analysis of different forecasting methods on the Exchange rates data, discussing their respective merits and drawbacks.

1 Introduction

Time series forecasting is the process of predicting future behavior using past signal correlations. Over the past decade, it has emerged as a significant area of study for both practitioners in the industry and academic researchers.

Numerous algorithms have been examined using variety of techniques, such as machine learning and deep learning methods (LSTM, RNN). [1] [2]

Nowadays, this learning problem is widespread in real-world applications where different observations are collected sequentially, such as medical data, electricity consumption temperatures or stock prices, hence the need to anticipate future changes, such as a fall in the exchange rate or a sudden rise in temperature, in order to remedy these perturbations and envisage appropriate solutions.

This project is divided into two parts:

- The first part is to provide a solution to the univariate forecasting problem of predicting the next 100 values of the OT variable for the ETTh1 dataset.
- The second part is aimed at finding a solution to the multivariate forecasting problem of predicting the next 100 values of the '6' variable of the exchange rate dataset, while dealing with the missing values.

2 Oil Temperature forecast

The first step is to train a predictor to accurately predict the next 100 values, which is a simple univariate prediction problem. To do this, we'll need to preprocess the data and split the dataset into a subset of subsequences of the form (input, target) for training. The training data consists of a univariate

time series with 1 feature. For this method, we will explain the techniques and their results.

We will also give a brief overview of the preprocessing steps we performed.

2.1 Prophet approach

Prophet is a predictive tool developed by Facebook in 2017. It's designed to analyze time series that show patterns on different time scales, such as annually, weekly and daily.

It also has advanced capabilities to handle the holiday effect, which is a very important factor in many practical cases. [2]

Here's a high-level overview of how Prophet works: The Trend Model in Prophet attempts to fit the data into a piecewise linear or logistic growth curve trend. This allows the model to handle non-linear growths by introducing change points at which the rate is allowed to change.

The Seasonality Model in Prophet is modelled using Fourier series. By default, Prophet will model seasonality with a period of one week and one year, but this can be set to other periods.

Holiday Effect: The user may provide Prophet with a list of dates that represent holidays or special events. Prophet will include these dates as additional regressors in its model.

Error Model: The error term is assumed to be normally distributed. The data is divided into a training set (80% of the data) and a testing set (the remaining 20%).

A Prophet model is created and fitted to the training data. The script then adjusts the predicted values by fitting a quadratic function to the data and generating new y values using this function. These new y values are then assigned back to the 'yhat' column in the forecast DataFrame. The script then saves the adjusted predictions and the original predictions to CSV files.

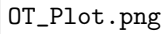
The image is a large, empty rectangular box with a thin black border, intended for a plot or figure.

Figure 1: Overview of the OT dataset distribution

2.2 Related works

The Prophet model is a statistical model. It employs past data to infer future trends, yet it lacks an understanding of the underlying process that generates the data. Consequently, its predictions are probabilistic, encompassing uncertainty intervals.

The Prophet model is particularly useful for business forecasting tasks due to its ability to generate forecasts in a relatively short timeframe, and its simplicity does not require a high level of familiarity with time-series data.

There are numerous related works that have explored different aspects of this technique. Taylor SJ and Letham B. [3] employed a decomposable time series model comprising three main components: trend, seasonality, and holidays. Meng, J, Yang, X, Yang, C and Liu, Y [7] conducted a deeply comparative experiment on the LSTM and the Prophet model on pharmaceutical sales prediction.

This is a brief summary of the research conducted on this model and the time series classification. It is an suitable introduction for anyone interested in this field of study.

2.3 Dataset and Pre-processing

On a selectionne las 15 monnaies les plus echangees aujourd'hui et on a extrait les differentes donnees historiques de chaque monnaie via Yahoo finance, comme ces differentes monnaies n'ont pas ete creees en meme temps on a pris un interval commun entre ces differentes monnaies qui est (23/12/2020- 24/04/2024)

The 'OT' values are distributed between -4.08 and 46.007, with a mean of 13.352. The majority of the dataset is represented by the value of 18.15, which is the 75th percentile. on a telecharge ces donnees dans le format csv et on a nettoye les donnees pour ne pas avoir de valeurs manquantes et on a creees une base de donnees relationnelle en utilisant MySql et Xampp.

ainsi on a obtenu 15 tables sql avec chacune une clef primaire qui est la date comme ca on pourra avoir une clef etrangere dans les autres monnaies et pour chaque table on a 4 colonnes open, close, high, low, volume. on appelle ca le format (OHLC+) (Citer la feuille de Daniel). Lutilisation d'un bdd SQL nous permettra d'avoir des donnees bien organisees qu'on peut manipuler en utilisant des requetes sql dans un programme python avec mysql.connector.

On peut essayer quelques requetes pour visualiser notre bdd sql.

(Ajouter des images du resultat de ces requetes)

Ensuite on ecrit des requetes sql pour extraire les differents donnees de notre bdd et de les mettre sous le format dataframe de pandas pour pouvoir les ajouter a notre modele LSTM qu'on va utiliser pour predire les valeurs futures de chaque monnaie.

on va ainsi essayer un LSTM avec un model Keras entraine sur 100 epochs et on repete le meme entrainement sur les autres monnaies. (On va aussi envisager d'essayer une approche multivariable des differentes monnaies et on va evidement calculer la correlation entre ces differentes monnaies)

On affiche les plots precedents.

On calcule la moyenne absolue d'erreur MAE et on essaye avec plusieurs Epochs et plusieurs monnaies.

on peut egalement envisager deprevoir les futures valeurs d'une monnaie en utilisant un modele entraine sur une autre monnaie (Par exemple on peut predire les valeurs futures du Bitcoin en utilisant un modele entraine sur la monnaie ETH)

on repete cette meme procedure sur l'ensemble des monnaies.

On compare le graph des monnaies depuis leurs creation jusqu'a aujourd'hui pour voir si les tendances de chaque monnaie est pareille durant toute sa duree de vie (hausse puis baisse puis hausse puis stabilisation)

On peut calculer le score F1 et les differents facteurs entre les monnaies pour avoir une idee de la similarite entre les differentes tables.

afficher les standard deviation et les pourcentilles.

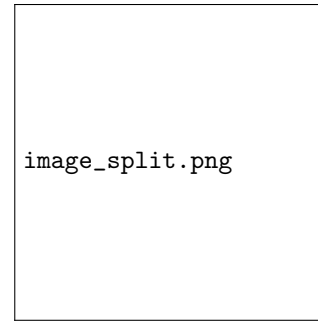


Figure 2: Distribution of the Train/Test sets

2.4 Training

On peut jouer sur les parametres d'entrainement du modele Keras et le LSTM comme le nombre d'epochs et la largeur du Batch, on essaiera le meme modele sur des splits 80/20, 60/40, 90/10. Si les predictions sont superieures aux valeurs precedentes alors SELL, s'ils sont inferieures alors BUY, si elles sont legerement differentes (on donne une petite marge d'erreur) alors on HOLD. Ainsi on peut garder notre benefice sans avoir a acheter ou a vendre.

On essayera nos estimations sur 3, 5 et 7 jours pour ne pas perdre en precision de predictions.

on repete ca avec les autres monnaies et on compare le graph des predictions et on peut tirer des conclusions (Par rapport a la similitude entre les differents plots ou entre les evenements reels comme la guerre d'Ukraine ou Israel).

We can also obtain plots of the different patterns (Figure3) with the command

```
model.plot_components(forecast)
```

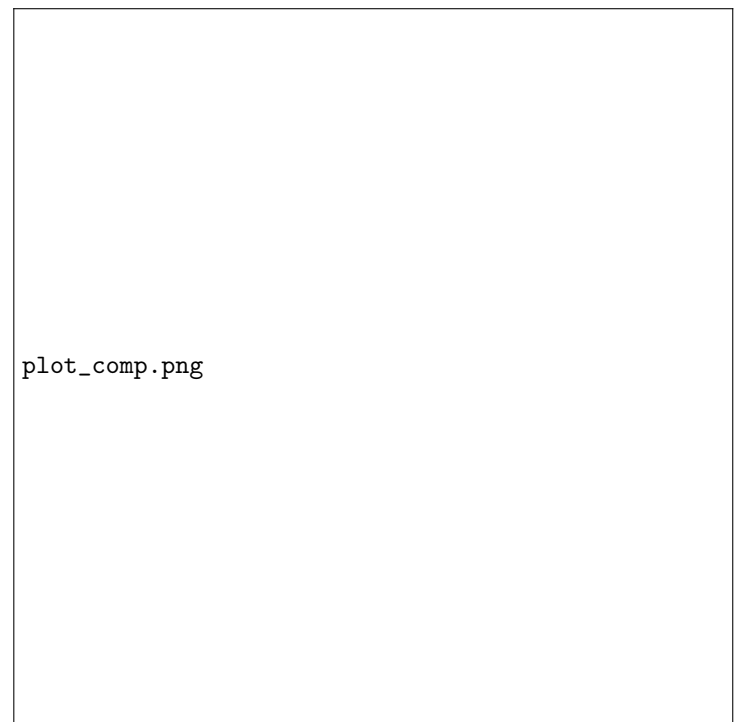


Figure 3: Plot of components

After running the model's training on the train set we can now define the desired time interval and the frequency (Hourly in our example) for predicting the 100 future values (until 28/06/2018 -17420 time series in total). To do this we use the command:

```
model.make_future_dataframe(periods = 100, freq = 'H')
```

After extracting the predicted values we can plot the whole result with matplotlib (Figure 4).



Figure 4: Plot of the future predictions

img_head.png

Figure 5: Example of the Tendencies

2.5 Indicateurs

On doit calculer des indicateurs utilises en commerce et en finance pour prevoir les monnaies en plus des OHLC+, on va les decire indic1 indic2 indic3 indic4

To evaluate this solution's performance we used the MAE (Mean Absolute Error) metric:

$$MAE = \frac{1}{N} \sum_{i=0}^N |Y_i - f(X_i)|$$

The scores obtained on the train/test sets and Kaggle sub-

mission were $\text{mae}_{\text{tt}}=0.7122$ and $\text{mae}_f=0.8912$, respectively.

2.6 Problemes

Le probleme avec le marche des crypto monnaies c'est qu'il est volatile et change quasi instantanement donc le defi majeur serait d'avoir des donnees chaque minute ce qui va generer un anorme taux d'informations qui va rapidement devenir hors de control, ainsi dans une semaine pour une seule monnaie on aura (10 080 donnees en collectant ces donnees pendant 5 ans on depassera les 2 419 200 donnees ce qui equivaut a 200Mb par fichier) In addition, the dataset and data pre-processing steps were discussed. Finally, the training parameters of the model were presented.

This section provides a comprehensive overview of the methods used for OT values forecast and the steps taken to prepare the data for training the Prophet model.

Although this model does not produce outstanding results, it is easy to deploy and can be implemented quickly.

3 Exchange Rate prediction

In this second part of the project, we're going to provide a solution to tackle the forecasting problem, to forecast the next 100 values of the '6' variable for the exchange rate dataset while handling the missing values.

3.1 Related works

There are numerous methodologies for performing multivariate forecasting. Two prominent examples are the RNN (Recurrent Neural Network) and the LSTM (Long Short-Term Memory) approach, as illustrated in the A. Sagheer, M. Kotb article [8]. In this study, the authors employed deep LSTM recurrent networks to predict petroleum production. Additionally, the combination of Prophet or Neural Prophet with a regressor enables the prediction of variables based on numerous inputs, as demonstrated in the Oskar.T,1,H.Hewamalagec, Polina P, N.Laptev, C.Bergmeir study [9], which compared the two previously mentioned models. Furthermore, the article JAGANNATHAN.J, DIVYA, C [10] presents an example of multivariate Prophet being used to predict climate change.

3.2 Dataset and Pre-processing

In order to achieve our objective, we have constructed a dataset comprising 7489 daily currency exchange rate time series, each of which represents a different data point (columns [0-6] and OT). However, 25% of the data is missing (33283 missing values out of the original 59912).

The following table shows the head of the original table.

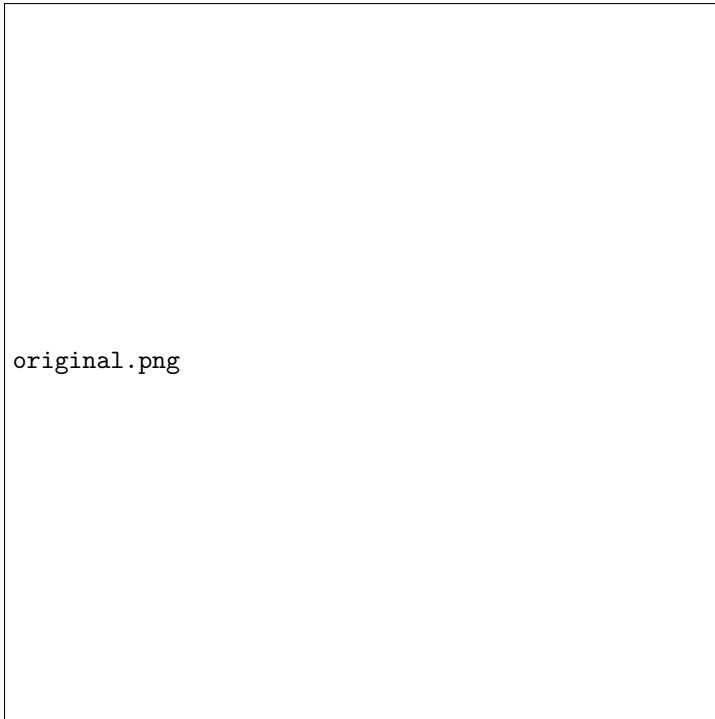


Figure 6: Exchange Rate with missing values

In order to overcome the problem of missing values, three methods of completion were tried:

- Mean Completion: Each of the missing values are replaced with the mean of the column `df.fillna(df.mean())`
- Median Completion: The missing values are replaced with the median of the column. `df.fillna(df.median())`
- Interpolation Completion: The missing values are filled based on a second order (quadratic) polynomial interpolation. `df.interpolate(method='polynomial', order=2).fillna(method='ffill')`, we used the second degree polynomial to avoid data overfitting.

We obtained the three following results (Figure 7, 8 and 9)

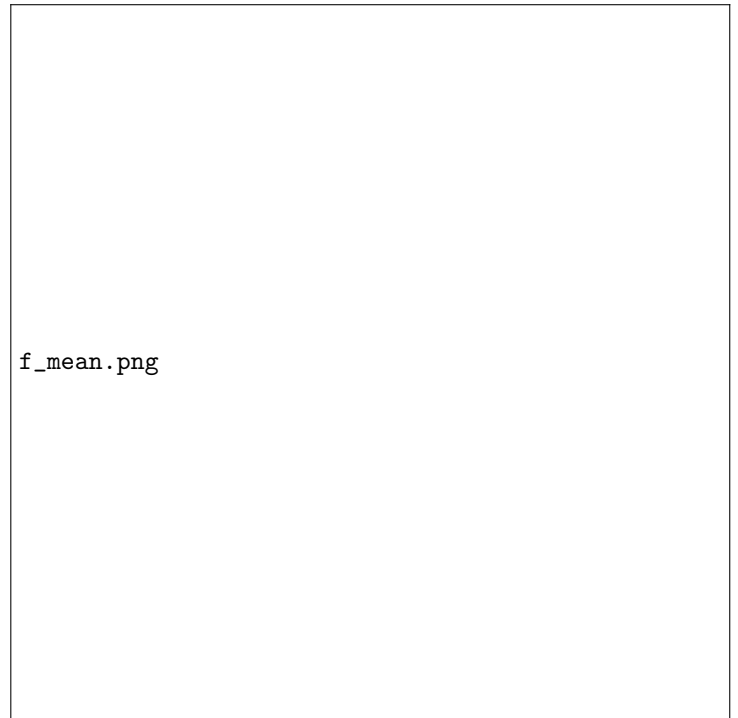
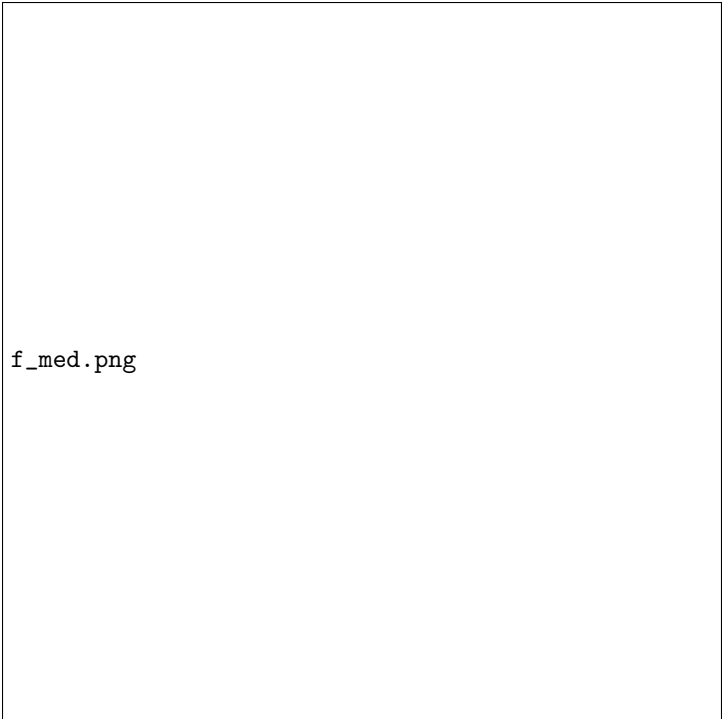
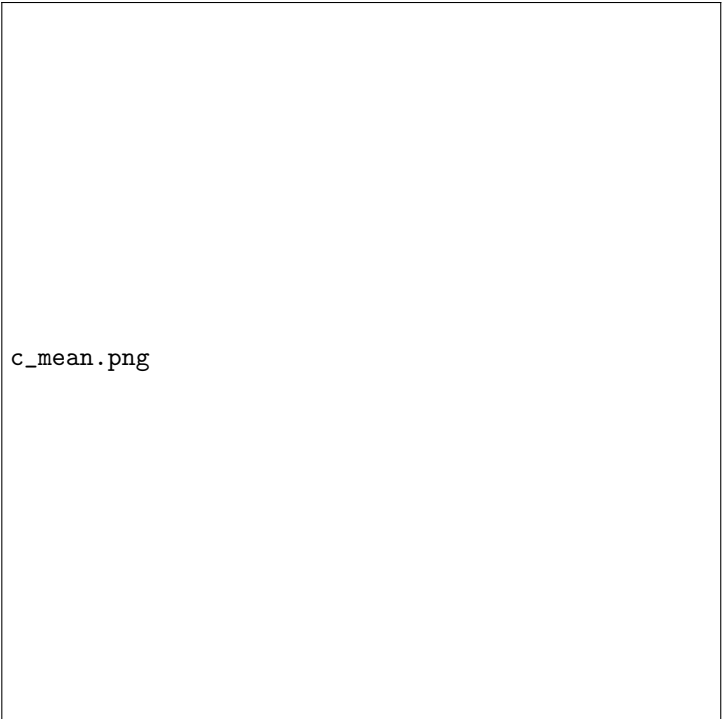


Figure 7: Mean Completion

A large rectangular box representing the median completion of the dataset. The text 'f_med.png' is located in the bottom-left corner.

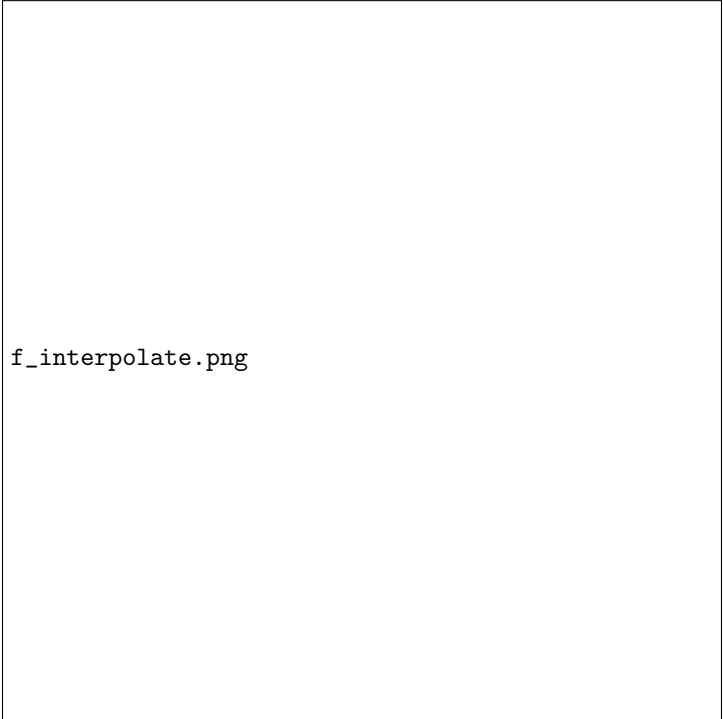
f_med.png

Figure 8: Median Completion

A large rectangular box representing the data distribution mean completion. The text 'c_mean.png' is located in the bottom-left corner.

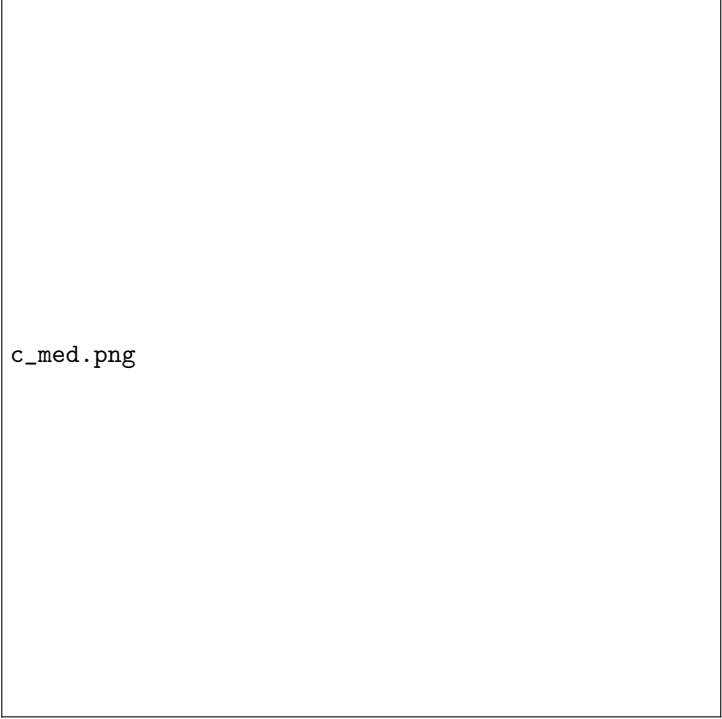
c_mean.png

Figure 10: Data distribution mean completion

A large rectangular box representing the interpolation completion of the dataset. The text 'f_interpolate.png' is located in the bottom-left corner.

f_interpolate.png

Figure 9: Interpolation Completion

A large rectangular box representing the data distribution median completion. The text 'c_med.png' is located in the bottom-left corner.

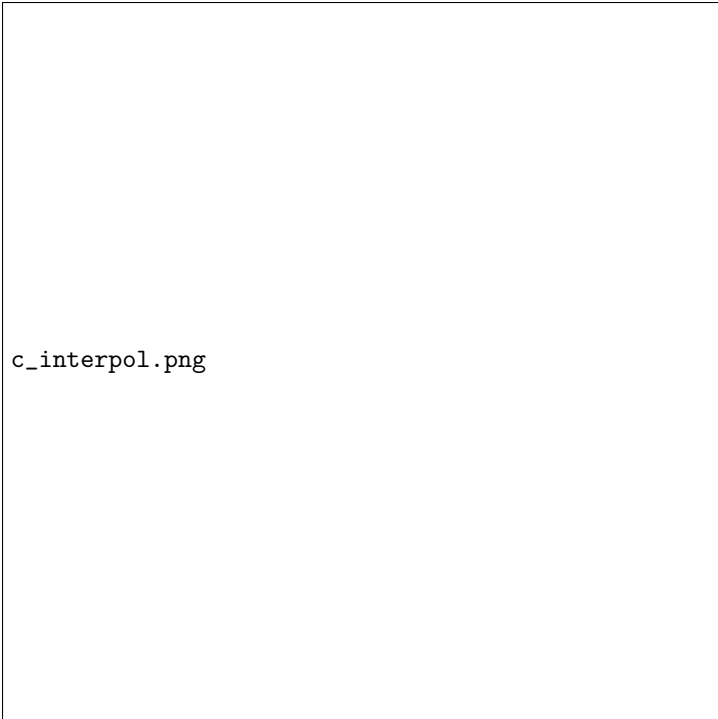
c_med.png

Figure 11: Data distribution median completion

3.3 Experiments

In order to address the problem, three different approaches were tested:

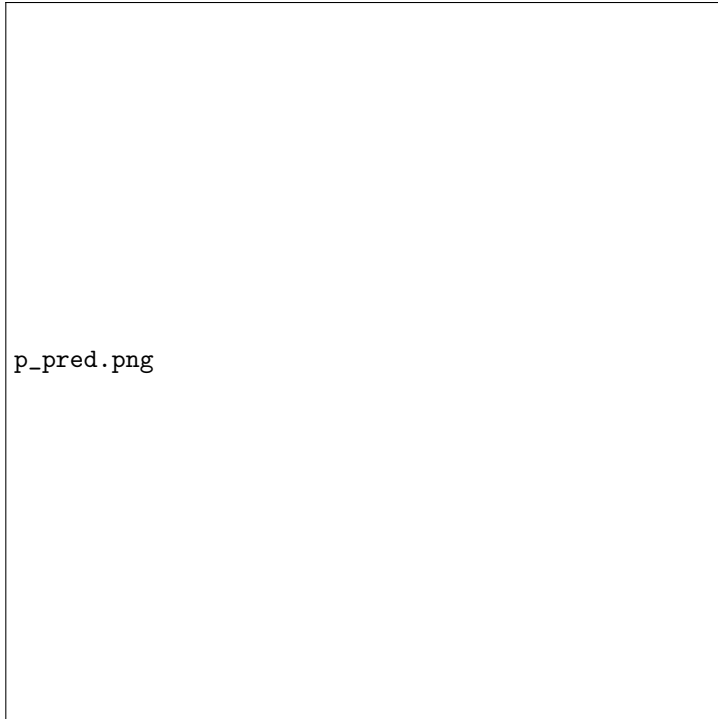
1. Univariate Prophet : We Applied the Step A method on this problem, and we tested it on the three completed version of the dataset(Figure 10, 11, 12)



c_interpol.png


Figure 12: Data distribution interpolation completion

The two columns, previously designated as 'date' and '6', were renamed as 'ds' and 'y', respectively. The dataset was then divided into two distinct sets, 80% of which constituted the training set and 20% the test set. This resulted in a total of 5990 instances in the training set and 1498 instances in the test set. The training set was then fitted into the Prophet model, resulting in the prediction of the figure 13 and 14 .



p_pred.png

Figure 14: Predictions using interpolation completion and Prophet



img_pred.png

Figure 13: Predictions using univariate Prophet model

2. LSTM (Long Short-Term Memory): Once the dataset had been completed by the preceding techniques, we proceeded to analyse the '6' column.

In order to ensure that the values were within the appropriate range for LSTM models, we applied a MinMax scaler to manipulate them between 0 and 1. In order to avoid having imbalanced sets, we then split the data into 65% Train and 35% Test. Subsequently, the input was reshaped to a format of [samples, time steps, features], which is the requisite format for LSTM.

The Sequential model was then created using Keras, which is a type of Recurrent Neural Network (RNN) that employs Long Short-Term Memory (LSTM) units, which are a specialised form of RNN unit that is effective for dealing with long-term dependencies in sequence data.

We executed the command `model.summary()` to obtain our model's architecture (Figure 15).

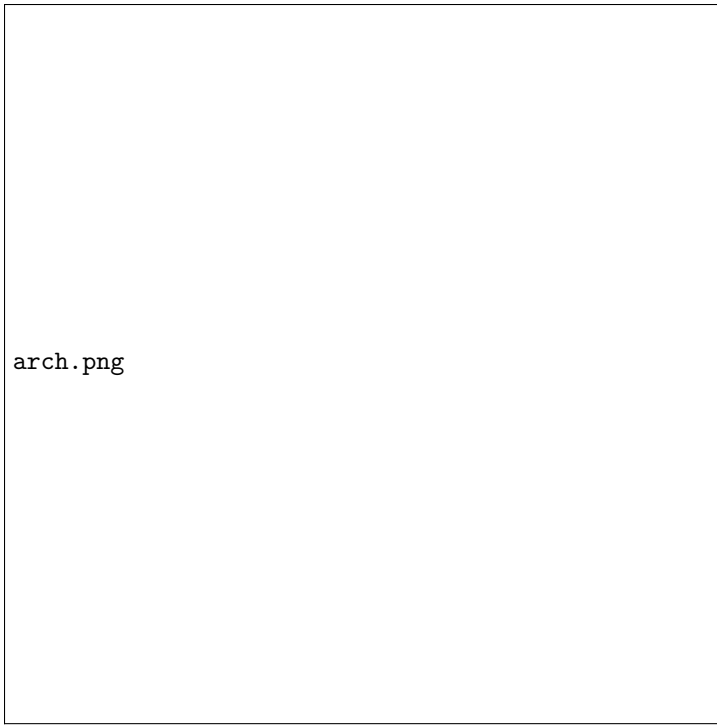


Figure 15: The Model's architecture

The data was then incorporated into the Keras model and training was initiated for 20 epochs with a batch size of 64.

The model's performance (Loss function [11]) was evaluated after each epoch.

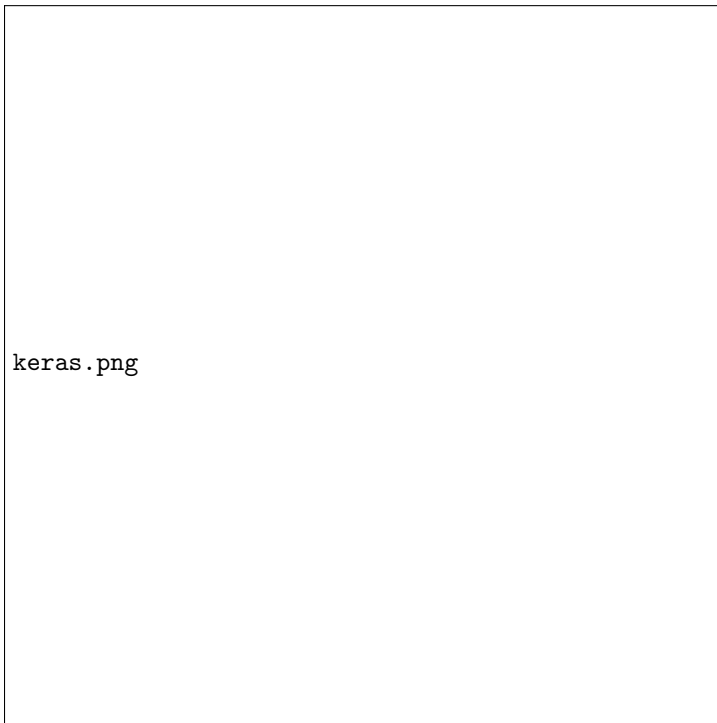


Figure 16: Keras model's performance after every epoch

Subsequently, the model's prediction (Blue) was displayed on the train (Orange) and test (Green) sets of

the 100 most recent values (Figure 17).

Finally, an inverse transform was applied in order to eliminate the data scaling, resulting in the plot of the predicted values (Figure 18).

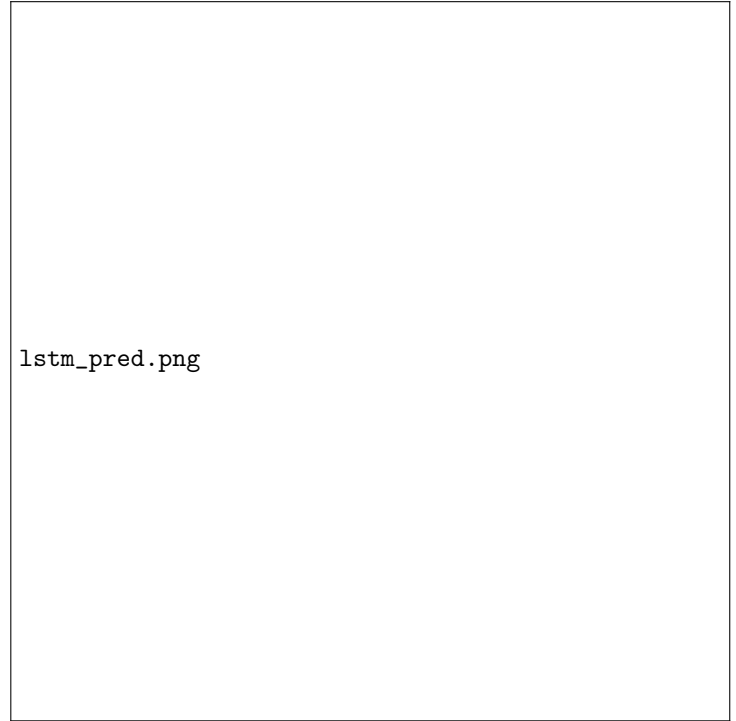


Figure 17: Model's prediction on train and test sets

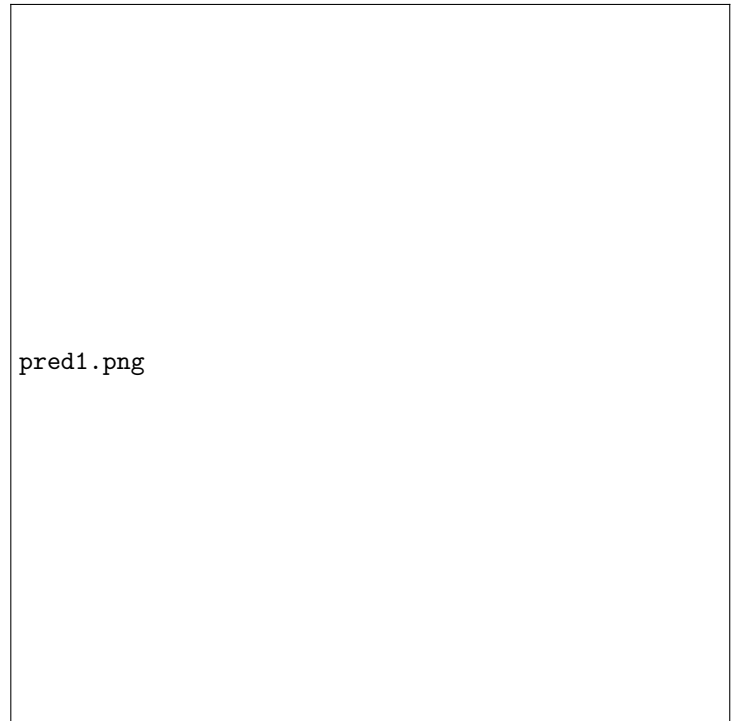


Figure 18: LSTM model's prediction

The same methodology was employed on the other completed datasets (mean and median), the results of which will be discussed later.

3. Multivariate Prophet : Now, we employed the previously cited prophet model in a multivariate approach. The dataset, which had been completed by the interpolation method, was imported and the start and end dates of the training and testing sets were defined. (Figure 19).

We plotted the different data distribution on the same figure to observe the relation between the variables visually (Figure 21).

Then we rename the date and the '6' column to 'ds' and 'y' and we select the relevant columns (Figure 20),

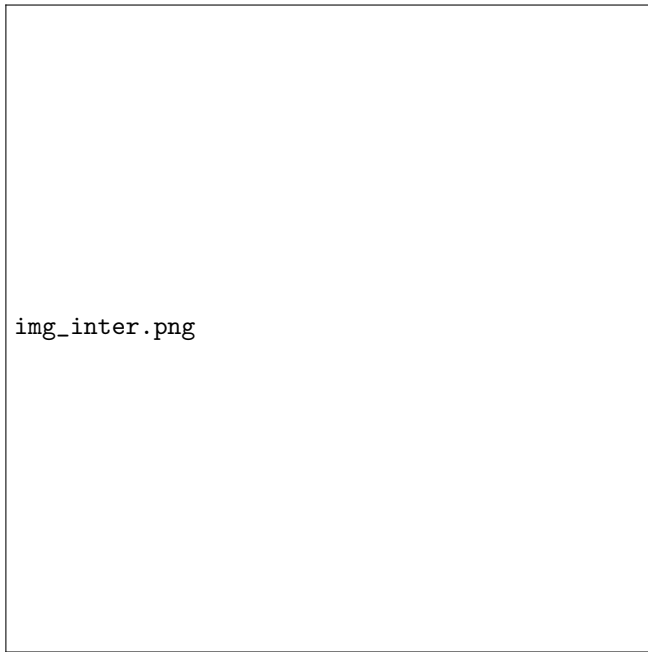


Figure 19: Head of completed dataset

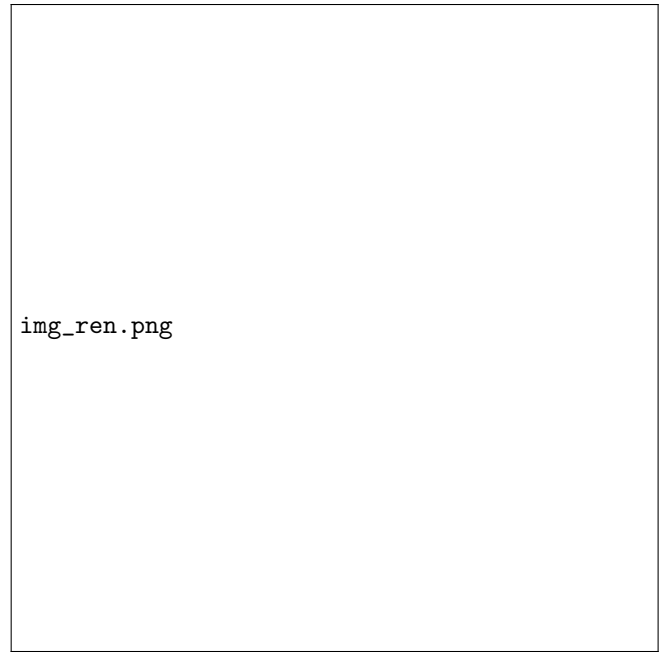


Figure 20: Relevant columns

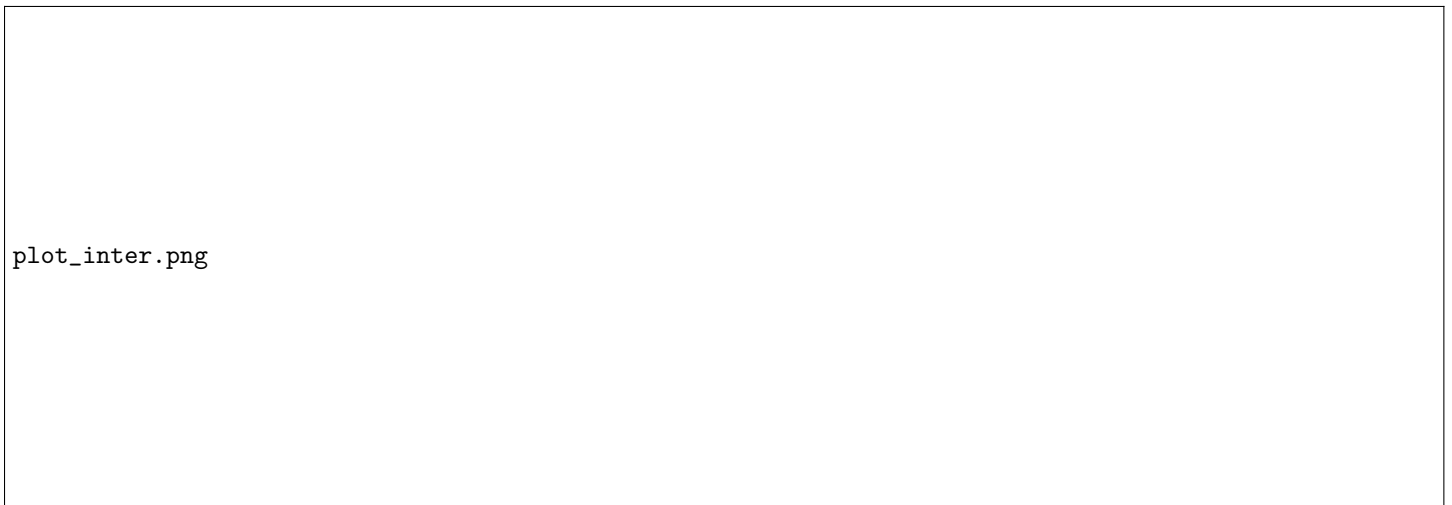


Figure 21: Dataset distribution

And to understand the relation between the different variables we calculated the correlation between each variable and the '6' variable (Figure 22)

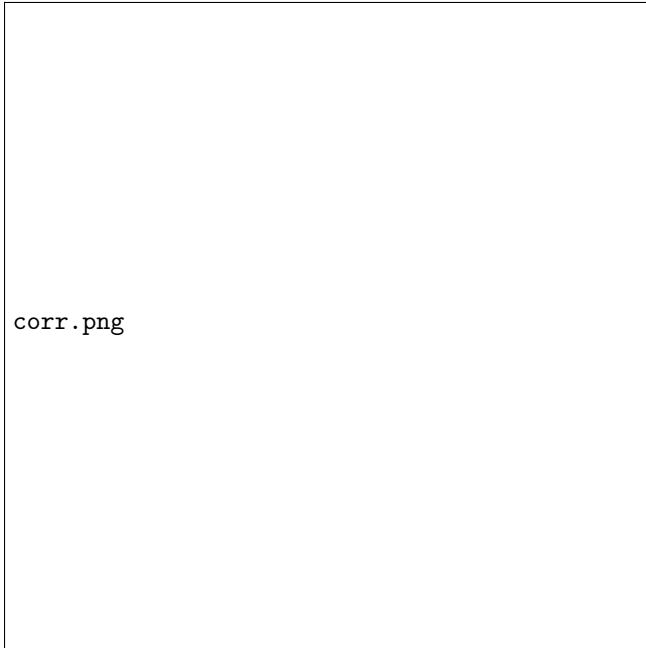


Figure 22: Correlation between variables

We notice that the best matching variables are the '3', 'OT' and '0'. These values are close to 1, indicating that there is a high correlation between the three columns and the '6' column.

We fitted the train/test sets into the Prophet model and we initiated the training process, then we defined a forecast time range of 100 and we plotted the forecast and the components chart.(Figure 23 and 24). In the forecast graphic, the light blue chart represents the uncertainty interval and the black dots represents the real values.

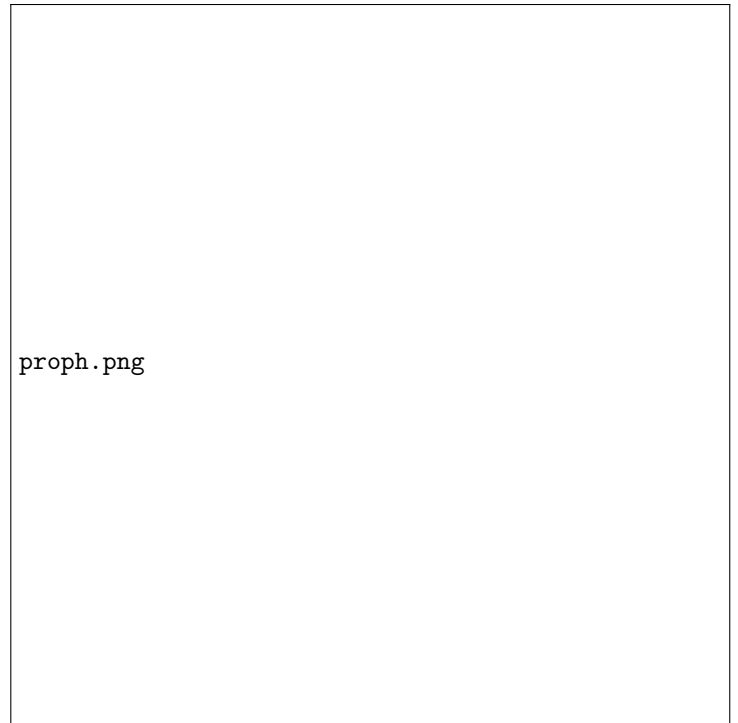


Figure 23: Forecast plot visualisation

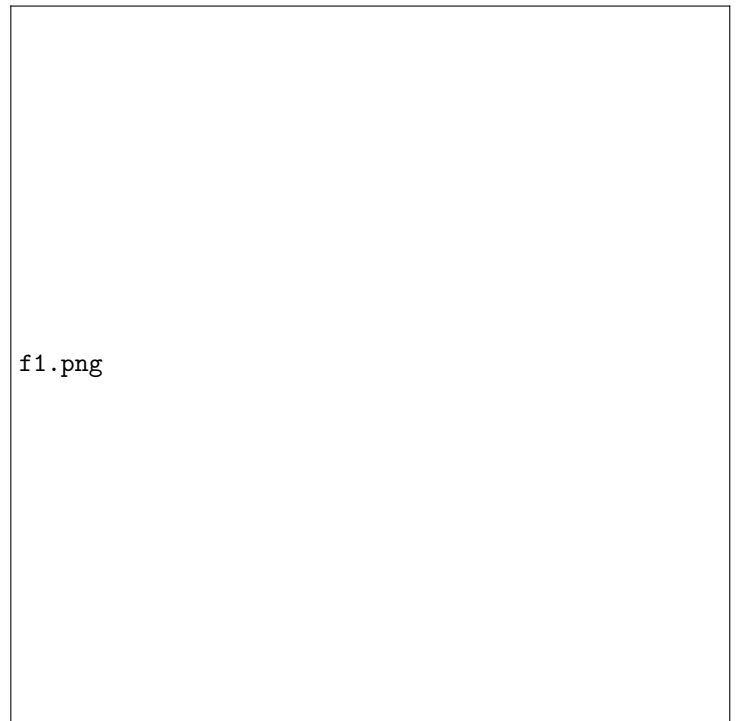


Figure 24: Forecasting components

By adding the weekly and the yearly seasonality to our initial baseline model, we displayed the 100 future predicted values and we plotted a second time the components (Figure 25).

f2.png

Figure 25: Future forecasting components

We observe that the yearly and the weekly tendencies didn't change very much from our baseline model components. However, the trend exhibited a notable change when the future prediction period was extended.

Finally, in order to enhance the predictive capabilities of our model, we incorporated an additional regressor comprising the columns '0', '3' and 'OT'. We standardised the input values within the regressor and tested the model's predictions for the subsequent 100 values on each case.

3.4 Results

To summarize the previous experiments, we have collected all the results in the table .

It is worth noting that the use of other correlated variables enhances the accuracy of our prediction, but we should also be careful to not use too many variables to avoid every column's estimation error (More columns=more error per column).

In our case, the optimal combination was the second-degree interpolation completion of the initial dataset and the use of the values of the column '3' to estimate the future values of the column '6'.

This approach yielded an mean absolute error of 0.0134.

3.5 Conclusion

In this section, we have presented the methodology that we have tested for forecasting the next values of the time series. In detail, we have explained the three models that were used. Additionally, we have provided a brief overview of related work and some of the research focused on LSTM and multivariate prophet methods in recent years. Furthermore, we have discussed the dataset and data pre-processing steps. Finally, we have provided details of the training parameters of the model. Furthermore, it should be noted that there are other methods of completing missing values that have not been discussed in this section(for example, the completion of missing data using association rules). The primary objective in this report was to identify a model that produces acceptable results and is both fast and easy to implement.

4 Summary

In conclusion, the two sections discuss different methods for forecasting future values from past time series data. The first section presents the univariate prophet method, which is a statistical technique that uses tendency and seasonality to predict certain patterns or motifs in the data. Although this method is fast and simple to implement, it can be less precise to determine the appropriate values if the dataset's size were voluminous.

Conversely, the second section presents three distinct models for performing multivariate prediction from exchange rates data. The third model outperforms the first and second models in precision and in execution time.

The results of the two sections are encouraging the establishment of a data monitoring system, whether to track the temperature of the oil or to monitor the exchange rate, with the objective of anticipating potential disruptions and implementing proactive repair measures.

Method	Approach	Completion	Mae	Columns used
Univariate	Prophet	Mean	0.3419	/
		Median	0.3483	/
		Interpolation	0.3292	/
	LSTM	Mean	0.0440	/
		Median	0.0633	/
		Interpolation	0.0357	/
Multivariate	Prophet	Median	0.1010	/
		Interpolation	0.0134	'3'
			0.0226	'0' and 'OT'
			0.0228	'OT'
			0.0223	'0' and '3'
			0.0197	'3' and 'OT'
			0.0240	'3' and '0' and 'OT'

Table 1: Summary of the previous experiments result

References

- [1] Makridakis, Spyros. Time series prediction: Forecasting the future and understanding the past. International Journal of Forecasting. 1994.
- [2] TOWARDS DATA SCIENCE. Time Series Analysis in Python: An Introduction. 2018. [consulted on 01/04/2024]. <https://towardsdatascience.com/time-series-analysis-in-python-an-introduction-70d5a5b1d52a>
- [3] MEDIUM. Prophet. 2020. Available at: <https://medium.com/@vinay1996/prophet-f729877e7e5a>
- [4] Taylor SJ, Letham B. 2017. Forecasting at scale. PeerJ Preprints 5:e3190v2 <https://doi.org/10.7287/peerj.preprints.3190v2>
- [5] P.Anand, M.Sharma, A.Sarolia, "Time Series Analysis for Predicting Anomaly using Prophet Models", AIJR Abstracts, pp. 70–70, Feb. 2024
- [6] Facebook, Inc. (2024). 'Uncertainty Intervals — Prophet 0.7 documentation'. Available at: https://facebook.github.io/prophet/docs/uncertainty_intervals.html (Accessed: 20 April 2024).
- [7] Meng, J., Yang, X., Yang, C., Liu, Y. (2021). Comparative Analysis of Prophet and LSTM Model. Journal of Physics: Conference Series, 1910(1): 12-59.
- [8] A.Sagheer, M.Kotb, Time series forecasting of petroleum production using deep LSTM recurrent networks, Neurocomputing, Volume 323, 2019, <https://doi.org/10.1016/j.neucom.2018.09.082>.
- [9] TRIEBE.O, HEWAMALAGE.H, PILYUGINA.P, et al. Neuralprophet: Explainable forecasting at scale. 2021. <https://doi.org/10.48550/arXiv.2111.15397>
- [10] JAGANNATHAN.J, DIVYA, C. Time Series Analyzation and Prediction of Climate using Enhanced Multivariate Prophet. International Journal of Engineering Trends and Technology, 2021. <https://ijettjournal.org/archive/ijett-v69i10p212>.
- [11] Ketkar, N. . Introduction to Keras. In: Deep Learning with Python. Apress, Berkeley, CA. 2017. https://doi.org/10.1007/978-1-4842-2766-4_7
- [12] Chih-Hung Wu, Chian-Huei Wun and Hung-Ju Chou, "Using association rules for completing missing data," Fourth International Conference on Hybrid Intelligent Systems, Japan, 2004, pp. 236-241, doi: 10.1109/ICHIS.2004.91.