1	Forecasting of Biogas Production Using Advanced Time Series
2	Algorithms
3	Nalluri Rishi Chaitanya Sri Prasad ¹ , Prabakaran Ganeshan ² , Karthik Rajendran ^{2*}
4	¹ Department of Computer Science and Engineering, School of Engineering and Sciences, SRM
5	University-AP, Amaravati, Andhra Pradesh 522240, India.
6	² Department of Environmental Science and Engineering, School of Engineering and Sciences, SRM
7	University-AP, Amaravati, Andhra Pradesh 522240, India.

8

*Corresponding author Email: rajendran.k@srmap.edu.in

Abstract 9

The anaerobic digestion (AD) process poses challenges in maintaining process 10 stability and time series-based prediction and forecasting due to the intricate nature 11 of the system. Process instability is a consequence of the unpredictability in the raw 12 13 material received at the facility, as well as temperature fluctuations and pH changes resulting from microbiological processes. Consequently, it is necessary to implement 14 15 constant monitoring and control measures for higher biogas production. The 16 challenges associated with anaerobic digestion (AD) systems can be effectively 17 addressed through the integration of advanced machine learning (ML) algorithms and 18 industry 4.0 systems within biogas facilities. This integration holds the potential to 19 enhance system efficiency and enable on-site control capabilities. Machine learning (ML) based solutions have the potential to enhance process performance in AD 20 21 facilities, leading to improved system operation and maintenance. The present study focuses on advanced ML techniques, specifically time series algorithms (ARIMA and 22 SARIMAX), have been employed to forecast the daily biogas production. These 23 algorithms are trained to discern critical process parameters and forecast daily biogas 24 production rates, measured in Liters. For forecasting 117 days of experimental data 25

26	used and identified ARIMA was best algorithm to forecast the daily production. This
27	algorithm excelled not only in predicting biogas production but also in forecasting
28	yield, resulting in a Root Mean Square Error (RMSE) of 3.26. Furthermore, a
29	comparison between the forecasted values of both ARIMA and SARIMAX was
30	conducted. The predictive ARIMA model underwent statistical validation with
31	unknown data, resulting in a P-value is >0.05.

32 Keywords:

Biogas; Anaerobic digestion; Modelling and prediction; Forecasting; Time series;Machine learning

35

36 Highlights

37	•	Forecast the biogas production with time series algorithms using ARIMA and
38		SARIMAX.
39	•	Feature Engineering and Hyperparameter tuning is done for better modelling.
40	•	Predicted and forecasted biogas production with an RMSE of 3.26 and 24.02
41		respectively.
42	•	By using this model, Biogas production can be forecasted for only the shorter
43		period.
44		

46 List of Tables

47	Table 1. Similar studies have been reported in literature	9
48	Table 2. Analytical follow-up of the co-digestion in the reactor	12
49	Table 3. Difference between the predicted yield and the actual yield	32

51 List of Figures

52	Figure 1 . Biogas plot with respect to days13
53	Figure 2. Overview of the Research Methodology Process
54	Figure 3. Statistics of the time series data
55	Figure 4. Seasonal decomposition of time series data24
56	Figure 5. ACF and PACF plots of ARIMA 26
57	Figure 6 . ACF and PACF plots of SARIMAX27
58	Figure 7. Plot diagnostics for ARIMA
59	Figure 8. Plot diagnostics for SARIMAX
60	Figure 9. Prediction of Biogas production in liters 30
61	Figure 10. Prediction of Biogas production in liters 32
62	Figure 11. Probability plot for the ARIMA and SARIMAX prediction with new data.
63	
64	Figure 12. Probability plot for the ARIMA and SARIMAX Forecasting of the biogas
65	production
66	

67 List of Equations

- **Equation 1.** AR(p) Model for ARIMA.
- **Equation 2.** MA(q) Model for ARIMA.
- **Equation 3.** Integrated (d) Model for ARIMA.
- **Equation 4.** Integrated difference (d=1) model.
- **Equation 5.** Integrated difference (d=2) model.
- **Equation 6.** ARIMA (p,d,q) model.
- **Equation 7.** SARIMA (p, d, q) (P, Q, D, s) model.
- **Equation 8.** Akaike information criterion (AIC).
- 76 Equation 9. Bayesian information criterion (BIC).
- **Equation 10.** Mean absolute error (MAE).
- **Equation 11.** Root Mean Square error (RMSE).
- 79 Equation 12. Additive model.
- 80 Equation 13. Multiplicative Model.

85 List of Abbreviations

AD	Anaerobic Digestion
SRT	Sludge Retention Time
VS	Volatile Solids
TS	Total Solids
BOD	Biochemical Oxygen Demand
COD	Chemical Oxygen Demand
SS	Total Suspended Solids
TN	Total Nitrogen
TP	Total Phosphorus
VFA	Volatile Fatty Acids
HRT	Hydraulic retention time
DNN	Deep Neural Networks
ML	Machine Learning
ARIMA	Autoregressive integrated moving average
AR	Autoregressive
MA	Moving Average
Ι	Integrated
SARIMAX	Seasonal Autoregressive integrated moving average
RMSE	Root Mean Square Error
NRMSE	Normalized Root Mean Square Error
MAE	Mean Absolute Error
MW	Megawatt

CAGR	Compound Annual Growth Rate
LSTM	Long Short-Term Memory
DA-LSTM	Dual-stage-attention based long short-term memory
DA-LSTM-VSN	Dual-stage-attention based long short-term memory
	integrated with variable selection networks.
RF	Random Forest
KNN	K-Nearest Neighbors
SVM	Support Vector Machine
ANN	Artificial Neural Network
Xgboost	Extreme Gradient Boosting
SELU	Scaled Exponential Linear Units
RMSprop	Root Mean Squared Propagation
AIC	Akaike information criterion
BIC	Bayesian information criterion
ADF	augmented Dickey-Fuller test
PACF	Partial Autocorrelation Function
ACF	Autocorrelation Function
NAN	Not a Number

91 **1. Introduction**

92 Utilizing the promising approach like anaerobic digestion (AD) for the manufacturers achieve the production of biogas, guaranteeing improved yields. The process of anaerobic digestion 93 involves microorganisms like acetogens, acidogens, and methanogens breaking down organic 94 95 matter such as animal manure, wastewater, and food waste, in the absence of oxygen (Nguyen et al., 2019). AD is efficiently designed and free from risks, ensuring that it structures its 96 97 processes in a manner that preserves human health. The outcomes of the anaerobic digestion process were, biogas (energy) and digestate (manure). The central primary focus revolves 98 99 around biogas production. Biogas consists of 25% - 45% carbon dioxide and 50% - 70% methane and some other traces gas. The usage of biogas finds extensive application in the 100 101 heating, electricity, and transportation sectors, serving as renewable fuels (Abanades et al., 102 2022).

103 In 2013, the worldwide recorded biogas production capacity reached approximately 14.173 104 MW, and this capacity increased to 18.505 MW by 2018, signifying a compound annual growth 105 rate (CAGR) of 30.69%. By 2022, the capacity further increases to 21.512 MW, demonstrating a production rate increase of 16.24% between 2018 and 2022. Over the span of 10 years, from 106 107 2013 to 2022, the potential for biogas production escalation is estimated to be around 51.84%. 108 This decade-long surge in biogas production exemplifies a substantial achievement in growth. 109 The potential of biogas is widely spread for the utilization across diverse sectors. The distribution of biogas percentages across sectors indicates allocations of 45% in electricity, 35% 110 111 in heating, and 10% in transport. In 2022, the global biogas market held a valuation of \$71.59 billion, and this value escalated to \$78.25 billion in 2023, showcasing a growth rate of 9.3% in 112 113 CAGR. The projection for global market growth in 2027 stands at an estimated \$102.7 billion, featuring a CAGR of 7.0%. 114

Table 1. Similar studies have been reported in literature.

Algorithm	Parameter	Best	Emphasis of	Model	Author(s)
used	considered	Algorithm	the study	Evaluation	
DA-LSTM- VSN, DA-LSTM, LSTM	Continuous: Qsludge, Qsludge, SRT, temperature. Discontinuous: VS/TS ratio, BOD, COD, SS, Total	VSN	Prediction of biogas production in anaerobic co- digestion of organic wastes	R ² = 0.76 NRMSE = 0.09	(Jeong et al., 2021)
	nitrogen (TN), and total phosphorus (TP)		using deep learning models		
RF, KNN, SVM, ANN, Xgboost	%TS, %VS, temperature, pH, alkalinity, VFA, CH ₄ and CO ₂ .	RF	Prediction of biogas production of industrial scale anaerobic digestion plant by machine learning algorithms	R ² = 0.92 RMSE=1405.8	(Yildirim and Ozkaya, 2023)
DNN	COD, BOD, TSS, PH, T, Total Nitrogen (TN), Total Phosphorus (TP)	DNN	DNN model development of biogas production from an anaerobic wastewater treatment plant using Bayesian hyperparameter optimization	R ² = 0.712	(Sadoune et al., 2023)
Extra Tree	TS%, VS, VFA mg/L, ALK mg/L, and VFA/ALK ratio	Extra Tree	Tree-Based Automated Machine Learning to Predict Biogas Production for Anaerobic Co- digestion of Organic Waste	R ² = 0.72	(Wang et al., 2021)
KNN, ARIMA, Decoder- ANN, Temporal Fusion Transformer	Raw sludge dry matter load, Sludge loss on ignition, Hydraulic retention time, Overnight stay, Raw sludge dry matter, Temperature, Sludge dry matter load, public holiday, Time index, Ambient temperature,	Temporal Fusion Transformer	Machine learning for quantile regression of biogas production rates in anaerobic digesters	RMSE=246.1	(Sappl et al., 2023)

Co-fermentation		
biowaste,		
Month, Sludge total,		
Sludge,		
Raw sludge,		
pH value,		
Sludge dry matter,		
Weekday		

117

118 In literature, the authors showed the comparison of biogas predictions using various 119 algorithms are represented in Table 1. (Jeong et al., 2021) They assembled time series data to 120 facilitate the training of diverse algorithms in deep learning. For both continuous and 121 discontinuous time series data, they utilized DA-LSTM-VSN as the optimal model for biogas prediction, achieving a test R² of 0.76 and a normalized root mean square error (NRMSE) of 122 123 0.68. Trained DA-LSTM with continuous input features, resulting in a test NRMSE of 0.10, 124 while LSTM, which was also trained with continuous input features, resulted in a test NRMSE 125 of 0.14. This model exclusively focused on biogas prediction and did not venture into future forecasting. (Yildirim and Ozkaya, 2023) Named a total of 8 input features to predict biogas 126 127 production, utilizing 5 machine learning algorithms. Among these, Random Forest (RF) stood out as the top performer, attaining an R² of 0.92 and an RMSE of 1405.84. Although the R² 128 129 value showcases satisfactory performance, the notable error rate indicates instances of high model overfitting. Furthermore, the model restricted its application to singular sample 130 131 predictions, and it did not include provisions for forecasting. (Sadoune et al., 2023) Deep Neural Networks (DNN) to train the model using 7 parameters. To enhance data modelling 132 through scaling, utilized three different scaling techniques: Minmax scaler, Robustscaler, and 133 Standardscaler. Among these techniques, Robustscaler emerged as the most effective for the 134 135 neural network, resulting in an R² of 0.712, considering a dense layer of 1 and a dropout rate of 0.144, activation function (SELU), and optimizer (RMSprop). The model was limited to 136 singular sample predictions without delving into forecasting. (Wang et al., 2021) Extra Tree 137 138 Regressor algorithm with 5 parameters as a tree-based model for biogas prediction. Artificial

139 Neural Network (ANN) modelled for result comparison with the Extra Tree algorithm. The 140 tree-based algorithm yielded an R² of 0.72, whereas the ANN algorithm resulted in an R² of 0.56. Their research mainly focus was solely on prediction, without incorporating biogas 141 forecasting into the study. (Sappl et al., 2023) Total of 4 prediction algorithms for biogas 142 forecasting, with one of them being a time series algorithm, namely ARIMA. All these 143 144 algorithms underwent consideration and were trained with the dataset. Among them, the Temporal Fusion Transformers algorithm yielded the best result, achieving an RMSE of 145 146 246.17, outperforming the other algorithms. On the other hand, ARIMA generated a 147 prediction result with an RMSE of 569.67, which is notably high compared to all the other 148 algorithms. It's important to note that ARIMA solely focused on prediction and didn't involve 149 forecasting biogas production daily. Currently, no published work has achieved the development of a time series model for forecasting biogas production with significantly low 150 151 RMSE and MAE. This is a crucial concern as the demand for biogas is rapidly increasing and 152 holds high significance for the future. Various sectors will actively utilize biogas in the upcoming days. 153

154 In this study focuses on biogas production, utilizing a dataset covering 117 experimental days. 155 Leveraging the advanced algorithm of time series within the realm of machine learning techniques, the primary focus is on forecasting and predicting model outcomes. Prediction 156 involves determining a solitary value for a given parameter or day, while forecasting entails 157 158 the continuous projection of outputs in accordance with the data sequence. Hyperparameter 159 tuning was conducted to select optimal parameters for AIC and BIC, crucial for model 160 training. Matplotlib library was utilized for clear graph visualization. Through the comparison of model evaluation, particularly RMSE and MAE, the model was chosen for both 161 prediction and data forecasting. If the model's performance falls short, it's retrained with 162 varying AIC and BIC values. The entire methodology of this paper is depicted in Figure 2. 163

164 2. Methodology

165 2.1 Data Collection and pre-processing

The data was collected from previously published research article from Scopus 166 database (Bautista Angeli et al., 2022). which was conducted experiment on a batch 167 mode to produce daily biogas production for a period of 117 days. This 168 comprehensive experiment incorporated several crucial operating and process 169 parameters, including VS consumed, initial pH, and final pH, Temperature, HRT 170 (hydraulic retention time), CH₄ mean production, , CH₄ content. To facilitate a 171 comprehensive analysis of biogas yield production, the experiment timeline was 172 segmented into three distinct periods based on these parameters. In the initial period, 173 first 12 days, an initial pH was 6.8 and a final pH of 5.5 were measured and a CH₄ 174 composition in produced biogas was measured at 44% (±6%). After the completion of 175 initial period, period 2 extended up to 63 days, continuously from 13 day. In this stage 176 an initial pH of 5.8, a final pH was 7.0, and a % CH₄ (% biogas) was 49% (±8%). The 177 178 subsequent phase, period 3, covered the duration from the end of period 2 to till 117th day, comprising around 54 days. During this final stage, the experiment adhered to 179 an initial pH of 7.0 and a final pH of 7.0, and a measured % CH₄ (% biogas) content 180 181 was 52% (±3%), as detailed in Table 2 (Bautista Angeli et al., 2022).

182

183

 Table 2. Analytical follow-up of the co-digestion in the reactor

VS CH4 mean VS consumed Initial Final consumed production (%) pН pН (mL/gVS)(%) Period 1 94 (±21) 94 (±1) 44 (±6) 5.5 6.8 Period 2 353 (±115) 83 (±7) 49 (±8) 5.8 7.0 52 (±3) Period 3 321 (±39) 83 (±7) 7.0 7.0

The dataset encompasses 117 days, facilitating a day-by-day examination of biogas 184 yield production data. We undertook data cleaning to address null values or missing 185 186 entries within the dataset, utilizing the front fill method to ensure a comprehensive 187 dataset for biogas yield analysis. This cleaned dataset was harnessed for the application of machine learning (ML) techniques, specifically time series analysis, to 188 forecast and predict the biogas production. Two distinct time series algorithms, 189 namely ARIMA and SARIMAX, were employed to predict or forecast biogas 190 production. To performance of these models can evaluate by using Root Mean Square 191 Error (RMSE) and Mean Absolute Error (MAE), for ensure the robust model 192 assessment. Figure 1 illustrates the relationship between biogas production with the 193 number of days. This visualization provides a clear overview of the trend over time. 194 In Figure 2 understanding of the entire methodology, in-depth for visually outlines 195 the comprehensive process employed in this study. 196





13

199 **2.2 Time Series Algorithms**

200 **2.2.1 ARIMA**

Autoregressive Integrated Moving Average (ARIMA) constitutes a generalized model 201 derived from Autoregressive Moving Average (ARMA) in the realm of time series 202 analysis. Application of algorithms are supposed to time-based data inputs 203 predominantly called for time series algorithms (Brockwell and Davis, 2016a), with 204 205 these being tailored to predict or forecast the patterns within the data. Often, such data exhibits non-stationary behavior, and is characterized by variations in statistical 206 properties. To enable the effective use of time series models for forecasting, it's 207 imperative to transform this data into a stationary form. Achieving stationarity 208 209 ensures that the mean average of data points remains consistent throughout time. The augmented Dickey-Fuller test (ADF) serves as a reliable method to assess data 210 stationarity. The p-value is equal to or less than 0.05, else not the data is deemed 211 stationary. Conversely, when data lacks stationarity, differencing becomes a vital step 212 213 to induce stationarity. This process is repeated iteratively until the data achieves the desired stationary state (Brockwell and Davis, 2016b). 214

The ARIMA framework contains three integral components that come into play: Autoregressive (AR), Integrated (I), and Moving Average (MA). The AR segment, denoted as "p" in Equation (1), embodies the autoregressive model. It signifies the number of observed lags in the model, commonly referred to as the lag order (Brockwell and Davis, 2016c).

220 AR(p) model:

221
$$y_t = c + \left(1 - \sum_{i=0}^p \varphi_i y_{t-i}\right) + \varepsilon_t \tag{1}$$

222 Where, φ_i = AR coefficients, y_t = Time series at time t, y_{t-i} = lagged values of time 223 series, ε_t = error time at time t, c = constant.

The Moving Average (MA) component, denoted as "q" in Equation (2), represents the rolling window size, also referred to as the order of the moving average. This parameter defines the span of observations considered for calculating the moving average at a specific time point (Brockwell and Davis, 2016c).

228 MA(q) model:

229
$$y_t = c + \left(1 - \sum_{i=0}^q \theta_i \varepsilon_{t-i}\right) + \varepsilon_t$$
(2)

230 Where, y_t = Time series at time t, θ_i = MA coefficients, ε_{t-i} = residual error at a lagged 231 period, ε_t = error time at time t, c = constant.

The selection of "p" and "q" values is derived from the Partial Autocorrelation Function 232 (PACF) and Autocorrelation Function (ACF) plots from the data. These plots encompass the 233 upper and lower boundaries. If the error bands cross these boundaries, they indicate time lags 234 that significantly impact the time series. These crossing points dictate the values for "p" and 235 "q". Specifically, "p" is determined from the PACF plots, and "q" is extracted from the ACF 236 plots. Within the ARIMA framework, the "I" component signifies integration, denoted as "d" 237 in Equation (3). This value represents the number of differencing operations applied to the 238 239 data to achieve stationarity. Each difference operation reduces the data's non-stationary behavior, progressively rendering it suitable for time series analysis (Brockwell and Davis, 240 2016b). 241

242
$$L^d y_t = y_t - y_{t-d}$$
 (3)

For the first differencing (d=1) of the data to make it stationary Equation (4),

244
$$y_t^1 = y_t - y_{t-1}$$
 (4)

245 In some cases, second differencing (d=2) also needed to make data to stationary

246 Equation (5)
$$y_t^2 = y_t^1 - y_{t-1}^1$$

247
$$= (y_t - y_{t-1}) - (y_{t-1} - y_{t-2})$$

248
$$= y_t - 2y_{t-1} + y_{t-2}$$
(5)

249 During this process, the application of differencing is pivotal to achieve data stationarity, enhancing the predictive capabilities for forecasting. ARIMA's effectiveness is notably 250 pronounced in short-term prediction scenarios, relying primarily on past data for its 251 forecasting abilities. The optimal selection of "p" and "q" values is determined through 252 253 statistical analysis, specifically utilizing metrics like the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC). These metrics gauge the model's fit by assessing 254 the goodness of fit relative to the complexity of the model. Lower AIC and BIC values signify 255 a better-fitting model. This approach, ARIMA constructs a predictive model capable of 256 accurate forecasting based on the learned data. The ARIMA equation, denoted as ARIMA (p, 257 d, q), is represented in Equation (6), encapsulating the interplay of autoregressive, 258 differencing, and moving average components (Brockwell and Davis, 2016c). 259

260 ARIMA (p, d, q) model:

261
$$\left(1 - \sum_{i=0}^{p} \varphi_i L^i\right) (1 - L)^d y_t = \left(1 - \sum_{i=0}^{q} \theta_i L^i\right) \varepsilon_t$$
(6)

262 L = lag operator, d = differencing term.

263 **2.2.2 SARIMAX**

Seasonal Autoregressive Integrated Moving Average with Exogenous Variables 264 (SARIMAX) extends from the foundation of Seasonal Autoregressive Integrated 265 266 Moving Average (SARIMA) model. The inclusion of "X" in SARIMAX signifies its 267 capability to handle exogenous variables. This algorithm is particularly effective for data displaying seasonality, making it a fitting choice for forecasting and prediction. 268 Seasonal ARIMAX, denoted as SARIMA (p, d, q) (P, D, Q, s) in Equation (7), 269 270 encompasses seasonal components in addition to the ARIMA parameters. The seasonal cycle length, represented as "s," typically corresponds to 12 months. The 271 values "p," "d," and "q" are derived from the ACF and PACF plots of the ARIMA model 272 (Brockwell and Davis, 2016c), indicated in Equations (1) & (2) & (3). Before modelling 273 data need to achieve stationarity through differencing. Since SARIMAX is the 274 algorithm with a seasonal cycle of 12 months requires a 12-month shift, differencing 275 is performed iteratively to ensure stationarity. The number of times differencing is 276 conducted is denoted as "D" in the SARIMAX model. The autoregressive model "P" in 277 SARIMAX is determined by assessing the seasonal PACF plot, while the moving 278 average "Q" is deduced from the seasonal ACF plot (Shumway and Stoffer, 2017). 279

281
$$\left(1 - \sum_{i=1}^{p} \varphi_{i} L^{i}\right) \left(1 - \sum_{i=1}^{P} \varphi_{i} L^{is}\right) (1 - L)^{d} (1 - L^{s})^{D} y_{t}$$

282 $= \left(1 + \sum_{i=1}^{q} \theta_{i} L^{i}\right) \left(1 + \sum_{i=1}^{Q} \Theta_{i} L^{is}\right) X_{t} + \varepsilon_{t}$ (7)

283 Φ_i = Seasonal AR coefficients, Θ_i = Seasonal MA coefficients, s= length of seasonal 284 period, X_t = exogenous variables.

285 2.3 Hyperparameter Optimization

286 Hyperparameter tuning is a technique for refining model performance and controlling the learning process to train the model efficiently. It aids to identifying the optimal parameters, 287 leading to enhance the model performance. In the context of time series algorithm used for 288 289 data modelling and forecasting, finding the best autoregressive (p, P) and moving average (q, Q) values, both for seasonal and non-seasonal components, holds utmost significance. 290 Optimized values of these parameters contribute to improved model fitting, enhancing 291 292 prediction and forecasting the outcomes. The selection process involves leveraging Akaike 293 Information Criterion (AIC) and Bayesian Information Criterion (BIC) techniques (Brockwell 294 and Davis, 2016a). These methods assist in identifying the most suitable AR (p, P) and MA (q, 295 Q) values for both ARIMA and SARIMAX models. The chosen values exhibit lower AIC and BIC scores, aligning with better model performance. AIC and BIC values, calculated as per 296 Equation (8) and (9), carry no fixed range. Instead, they are influenced by the specific AR and 297 MA values derived from ACF and PACF plots. AIC and BIC scores need to be substantially 298 299 lower compared to alternatives, indicating the ideal choice of parameters for optimal 300 modelling.

301 Akaike information criterion (AIC):

$$AIC = 2k - 2\ln(L) \tag{8}$$

303 k = number of parameters in the model, L = max value of likelihood function of model

304 Bayesian information criterion (BIC):

305
$$BIC = k * \ln(n) - 2\ln(L)$$
 (9)

306 n = number of observations.

307 **2.4 Model Evaluation**

Time series algorithms are utilized for both prediction and forecasting of the model. These models were effectively implemented and adapted to the dataset. The model performance was evaluated based on the metrics of Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) (James et al., 2013). MAE represents the average of the absolute differences between the true and predicted values. As the difference between the true and predicted values increases, the RMSE also increases, as demonstrated in Equation (10).

315
$$MAE = \frac{\sum_{i=1}^{n} |y_i - \hat{y}_i|}{n}$$
(10)

316 Where, y_i = Actual values, \hat{y}_i = Predicted values, n = Number of data points

Root Mean Square Error is employed to calculate the squared differences between the actual and predicted values of the output variable. The RMSE provides insight into the proximity of the actual data points to the predicted data points. A lower RMSE signifies a minimal difference between the true and predicted values, as illustrated in Equation (11).

322
$$RMSE = \frac{\sqrt{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}}{n}$$
(11)

323 y_i = Actual values, \hat{y}_i = Predicted values, n = Number of data points



Figure 2. Overview of the Research Methodology Process.

326 3. Results and Discussion

327 3.1 Data Statistics before Modelling

The statistical analysis of the dataset involves 117 days of daily biogas production 328 329 experiments. Throughout this period, the biogas produced ranged from a minimum of 1.86 to a maximum of 405.9 liters, with a mean production of 153.14 liters per day. 330 The experiments were divided into three distinct periods, each characterized by 331 332 specific parameters such as initial pH and final pH, which were detailed in Table 2. 333 To address missing values, a front-fill method was employed. Additional insights into the data's distribution and variability can be assembled from the provided box plot in 334 335 Figure 3.



Figure 3. Statistics of the time series data

339

336

337

338

341 **3.2 Seasonal Decomposition**

Seasonal decomposition involves utilizing statistical analysis to dissect time series 342 data into distinct components such as trend, seasonality, and residual patterns. These 343 components offer valuable insights into the various variations present within the data, 344 345 aiding in comprehensive analysis. Two primary models utilized for seasonal decomposition are Additive and Multiplicative (Cowpertwait and Metcalfe, 2009). 346 347 Additive models are employed when the magnitude of seasonal data remains 348 consistent throughout the entire period without notable fluctuations. This model aggregates the seasonal component, trend, and residual. The formulation of the 349 additive model is presented as Equation (12). 350

$$y_t = T_t + S_t + R_t \tag{12}$$

The Multiplicative model comes into play when the magnitude of seasonal data exhibits changes throughout the time span. This model is employed in situations where the fluctuations in seasonal patterns are not consistent. It involves calculating the product of seasonal components, trends, and residuals. The formulation of the multiplicative model is represented as Equation (13).

$$y_t = T_t * S_t * R_t \tag{13}$$

358 T_t = Trend component, S_t = Seasonal component, R_t = Residual component.

In seasonal decomposition, the trend component signifies the smoothed trajectory of the time series data. This component offers insights into whether the time series exhibits an upward, downward, or steady trend. On the other hand, the Seasonal Component is utilized to identify recurring patterns that manifest at regular intervals within the time series data. These patterns can be associated with specific days, months, or years, and the periodicity can be adjusted based on analytical requirements. The Residual Component captures the residual noise that remains once the trend and seasonal patterns have been removed from the data. This segment encompasses short-term fluctuations that cannot be explained by the underlying trend and seasonality components.

We collected 117 daily experimental observations for the present time series data 369 370 concerning biogas production. Utilizing the multiplicative model for seasonal 371 decomposition was appropriate (Cowpertwait and Metcalfe, 2009), given the varying 372 nature of the data over time. Examining the decomposed components, the observed 373 component presented a normal distribution of the data. The trend component depicted a discernible upward trajectory, indicating an increasing trend within the 374 time series. In terms of the seasonal component, recognizable patterns were evident, 375 376 showcasing periodicity in the data. The residual component, denoting the error, fell within a range of 0.5 to 1.5, suggesting relatively minor discrepancies. The seasonal 377 decompose plot can be observed in Figure 4, visually encapsulating the outlined 378 379 components and their behavior.



380

381

Figure 4. Seasonal decomposition of time series data.

382

383 **3.3 Time Series Analysis**

384 **3.3.1 ADF (Augmented Dickey-Fuller test)**

The Augmented Dickey-Fuller test (ADF) is utilized to determine whether the data 385 exhibits stationarity. Stationarity ensures a constant average mean across all data 386 points over time. In the ADF test, you calculate a p-value; if this value is less than or 387 equal to 0.05, the data is considered stationary. To achieve stationarity, the technique 388 of differencing is employed. In the context of predicting biogas using the ARIMA 389 model, the ADF test initially yielded a p-value of 0.57, indicating non-stationarity of 390 the data. The first difference was applied by shifting the data by 1. Subsequently, the 391 ADF test for this first-differenced data revealed an exceptionally low p-value of 9.99e-392 07, confirming its stationarity. Transitioning to the Seasonal ARIMA (SARIMAX) 393 model, which accounts for seasonality, a significant shift of 12 was required for 394

differencing due to its seasonal nature. Remarkably, after applying this seasonaldifference, the SARIMAX model exhibited stationarity, supported by a p-value of 0.04.

397 **3.3.2 ACF & PACF plots**

Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) serve 398 as essential statistical analysis tools for enhancing time series models in the quest for 399 improved prediction and forecasting accuracy. ACF plots play a pivotal role in 400 determining the Moving Average (MA) model components, while PACF plots are 401 402 instrumental in identifying the Autoregressive (AR) model components. These plots play a central role in the development of ARIMA and SARIMAX algorithms. Both 403 404 ACF and PACF plots are accompanied by upper and lower bounds, which form a 405 confidence interval. When the error bands within these plots intersect or cross over 406 this confidence interval, it signifies the statistical significance of that lag. ACF plots 407 encapsulate the correlation between a data point and its own past lags. These plots are particularly adept at unveiling patterns within the data. When ACF values exhibit a 408 gradual decrease, it indicates the presence of a persistent pattern in the data. PACF 409 plots, on the other hand, reveal the specific relationship between a data point and its 410 immediately preceding value, effectively capturing the direct influence of past 411 observations on the current point. The range of ACF and PACF values spans from -1 412 to 1. Where, -1 indicates a negative correlation, 0 signifies no correlation, and 1 413 represents a positive correlation (Shumway and Stoffer, 2017). 414

In the ARIMA (p, d, q) model, we posit an AR(p) value of 40 based on observations
from the PACF plot. For the MA(q) component, we designate a value of 10, informed
by insights gleaned from the ACF plot. As for the differencing aspect denoted by 'd',

a value of 1 is selected. This decision aligns with the need to attain data stationarity,
which was achieved through a single differencing operation, as illustrated in Figure
5.



422

Figure 5. ACF and PACF plots of ARIMA

423 To build an accurate and efficient time series forecasting model, we also utilized the SARIMAX approach. The non-seasonal plots of the Autocorrelation Function (ACF) 424 and Partial Autocorrelation Function (PACF) guided us in determining the non-425 seasonal orders of differencing, represented as (2, 1, 2). For the seasonal component, 426 we focused on the ACF and PACF plots after applying a seasonal shift of 12, 427 corresponding to a yearly cycle. This allowed us to extract the seasonal orders of the 428 SARIMAX model, which we established as (2, 1, 1, 12) Figure 6. For identifying the 429 430 optimal parameters for our model, we relied on two key criteria: the Akaike 431 Information Criterion (AIC) and the Bayesian Information Criterion (BIC). By comparing different parameter combinations using these metrics, we were able to 432 select the parameters that provided the best fit for the time series data. Overall, these 433 comprehensive analyses of the non-seasonal and seasonal plots of ACF and PACF, 434 coupled with the utilization of AIC and BIC, facilitated the development of a robust 435

time series forecasting model. This model, based on SARIMAX, ensures a precise and
efficient prediction of future values, thereby enhancing our ability to make informed
decisions based on the insights derived from the time series data (Brockwell and
Davis, 2016a).



441

440

Figure 6. ACF and PACF plots of SARIMAX

442 **3.3.3 Plot Diagnostics**

The diagnostic analysis of the plotted statistics involves four main types of statistical 443 graphs: Standardized residuals, Histogram plots, Normal Q-Q plot, and Correlogram. 444 To assess the quality of the model's fit, various visualizations are utilized. 445 Standardized residuals portray the discrepancies between actual values and predicted 446 447 values. Calculating standardized residuals involves dividing normal residuals by the standard deviation of total residuals. This produces a standardized residual plot, 448 illustrating the standardized differences in errors between actual and predicted 449 values. Histogram plots showcase the distribution of residuals. The graph displays 450 two curves: the normal distribution curve N (0,1) and the Kernel Density Estimate 451 (KDE) curve. If the lines exhibit minimal disparity, the model is considered well-fitted. 452 The KDE, a smoothed version of the histogram, generates a continuous representation 453

of data density. Particularly for ARIMA and SARIMAX models, a Gaussian kernel is 454 employed for KDE plotting. Peaks in the KDE graph pinpoint areas of higher data 455 456 density, while lower regions signify lower density. This visualization effectively 457 communicates the data distribution's central tendency and variance. The Normal Q-Q plot gauges data distribution against the normal distribution. When data points 458 align with the red reference line, it indicates normal residuals. Both the ARIMA 459 (Figure 7) and SARIMAX (Figure 8) models closely adhere to this red line, signifying 460 a nearly normal distribution of residuals. 461



462

463

Figure 7. Plot diagnostics for ARIMA.

The Correlogram plot generates the Autocorrelation Function (ACF) plot using residuals instead of data. This plot reveals whether error bands fall within the confidence interval, indicating if residuals conform to a normal distribution. Consistency within the interval signifies normal distribution adherence, while deviations suggest potential data omission by the model. Notably, the lines of the ARIMA model fall insignificantly within the confidence interval, implying a lack of significance. Conversely, SARIMAX exhibits noticeable differences Figure 8. When 471 comparing both algorithms, it becomes evident that the ARIMA model predicts data 472 more accurately than SARIMAX. Upon model evaluation, the ARIMA model 473 demonstrates an RMSE of 3.26, whereas the RMSE for SARIMAX is 24.02. The 474 superiority of the ARIMA model in this regard is apparent. These plot diagnostics 475 leverage a variety of graphical tools to observe model residuals in time series analysis.



476

477

Figure 8. Plot diagnostics for SARIMAX

478 **3.4 Comparison of Time Series algorithms**

479 **3.4.1 Prediction**

480 We used time series analysis to predict the model by fitting two algorithms, namely ARIMA and SARIMAX, to the data. Both algorithms were applied to the output for 481 prediction, specifically biogas production, with the p, d, q values determined from the 482 Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) plots. 483 After observing AIC and BIC values of 1402.4 and 1542 for the order (40, 1, 10), the 484 ARIMA model was trained and found to be better fitted to the data. In the case of 485 SARIMAX, the seasonal orders (P, D, Q) were interpreted from seasonal ACF and 486 PACF plots, resulting in an order of (2, 1, 2) and (2, 1, 1, 12). The AIC and BIC values 487

were 1245.3 and 1266.5, respectively, for this SARIMAX configuration, which was then
fitted and used for prediction. Both ARIMA and SARIMAX algorithms provided
predictions that were well-fitted to the actual biogas production, as shown in Figure
9. The ARIMA model predicted daily biogas production with Mean Absolute Error
(MAE) and Root Mean Squared Error (RMSE) values of 46.81 and 3.26 respectively.
On the other hand, SARIMAX achieved MAE and RMSE values of 60.81 and 24.02.
Considering the time series modelling, both ARIMA and SARIMAX algorithms

494 Considering the time series modeling, both AkhviA and SAkhviAA algorithms
495 exhibited very low prediction errors for biogas production. This underscores the
496 effectiveness of these algorithms in accurately forecasting biogas production based on
497 the time series data.



498

499

Figure 9. Prediction of Biogas production in liters

500

501 **3.4.2 Forecasting**

Forecasting involves predicting future unknown days by training the model usingpast data. This type of modelling allows us to anticipate whether outcomes will

increase or decrease in the future. After fitting the model with the time series 504 algorithms, namely ARIMA and SARIMAX, the focus shifts to the forecasting phase. 505 506 To execute the forecasting technique, we introduce unknown days representing a 507 future period. These days are initialized with "Not a Number" (NAN) values to facilitate predictions for unfamiliar data. These NAN values are then combined with 508 509 the existing dataset, creating an extended dataset. This extended dataset is subsequently utilized to train the time series algorithms, enabling them to forecast the 510 511 future data points. In our dataset, information is available for 117 days, while the final dataset incorporates new days to facilitate forecasting. To define the range for 512 forecasting, we designate start and end dates as day 77 and day 150. This range 513 signifies that forecasting spans from the 77th day to the 150th day, though this range 514 can be customized to suit our requirements. Both ARIMA and SARIMAX algorithms 515 perform well and effectively forecast biogas production up to the 150th day. The 516 forecasting results of ARIMA and SARIMAX are visually represented in Figure 10. 517 This visual representation provides insights into the predicted trends and variations 518 in biogas production over the forecasting period. 519



Figure 10. Prediction of Biogas production in liters

522 **3.4.3 Testing model against new data**

To evaluate the model using unfamiliar data, we examined five new data points 523 corresponding to days when biogas experiments were conducted. These data points 524 525 were utilized to assess the variance between the biogas production predicted by the time series model and the actual biogas generated. Predictions using both the ARIMA 526 and SARIMAX algorithms, which yielded RMSE values of 3.26 and 24.02, respectively. 527 When we fed this data into the model for prediction, it yielded a significantly lower 528 error in comparison to the outcomes generated by the time series algorithm. All these 529 values are presented in Table 3 for reference. To evaluate the forecasting results, both 530 ARIMA and SARIMAX algorithms were employed to predict biogas production from 531 the 77th day up to the 150th day (as depicted in Figure 10). These forecasts were made 532 for the month. 533

534

521

Table 3. Difference between the predicted yield and the actual yield

535

Time Series	Day	Observed Biogas	Forecasted	Deviation (%)
Algorithm	-	_	Biogas	
ARIMA	5	71.96	61.35	10.61
	13	83.03	89.03	-6
	19	77.49	63.54	13.95
	25	55.35	62.49	-7.14
	35	162.36	118.6	43.76
SARIMAX	5	71.96	84.59	-12.63
	13	83.03	57.61	25.42
	19	77.49	128.4	-50.91
	25	55.35	41.86	13.49
	35	162.36	116.59	45.77

537 **3.4.4 Statistical Analysis**

This technique is utilized to test whether a statistically significant difference exists 538 between the two models, or it can also be employed to assess the statistical difference 539 between the actual and predicted values. This testing process is essential because, for 540 541 a model to align with the null hypothesis, the p-value must be greater than 0.05. This implies that the newly predicted values fall within the 95% confidence interval of the 542 543 data distribution, signifying no significant difference between the actual and 544 predicted values. If the p-value is less than 0.05, the null hypothesis is rejected in Favor of the alternative hypothesis. In such cases, the model is deemed to exhibit a 545 significant difference between the predicted and actual values (James et al., 2023). 546 Based on the data presented in Table 3, we extracted the disparity between the actual 547 and predicted models to execute the p-test. The outcome revealed a p-value exceeding 548 0.05, signifying that both models demonstrate an insignificant difference in their 549 predictive capabilities. This observation is consistent with the model comparison 550 shown in Figure 11. 551



Figure 11. Probability plot for the ARIMA and SARIMAX prediction with new data.

For the forecasting model, we calculated the discrepancy between the ARIMA and SARIMAX predictions from the 77th day to the 150th day. Subsequently, a paired ttest was conducted between these two sets of predictions. The obtained p-value was greater than 0.05, indicating that the forecasting model similarly suggests an absence of significant difference between these two models, even in the context of biogas production forecasting. This outcome is illustrated in Figure 12.



560

Figure 12. Probability plot for the ARIMA and SARIMAX Forecasting of
 the biogas production.

563 **4. Conclusion**

The culmination of this paper is the forecasting of biogas production through the utilization of time series algorithms, specifically ARIMA and SARIMAX. The dataset comprises 117 days of daily biogas production experimentation, organized into three distinct periods. Throughout each period, key parameters such as CH₄ mean production, VS consumed, CH₄ (% biogas), initial pH, and final pH are meticulously upheld, as detailed in Table 2. To address missing data for certain days, feature engineering techniques are employed. Additionally, hyperparameter tuning is

conducted to determine optimal values for ARIMA (p, d, q) and SARIMAX (p, d, q) 571 (P, D, Q, s) parameters. This process is guided by the AIC and BIC values to ensure a 572 573 well-fitted model for predictive and forecasting tasks. The identification of suitable 574 parameter values is facilitated by examining the plots of the autocorrelation function 575 (ACF) and partial autocorrelation function (PACF). Consequently, the dataset is fitted to the ARIMA and SARIMAX algorithms. The application of these models results in 576 the prediction and forecasting of biogas production for future unknown days. The 577 achieved RMSE values are 3.26 for ARIMA and 24.02 for SARIMAX, respectively. 578 Subsequent statistical analysis of the new data confirms p-values greater than 0.05 for 579 both prediction and forecasting, signifying the absence of significant differences 580 between the two methods. Thus, these models are deemed viable for biogas 581 production forecasting. 582

584 **5. References**

585	Abanades, S., Abbaspour, H., Ahmadi, A., Das, B., Ehyaei, M.A., Esmaeilion, F., El
586	Haj Assad, M., Hajilounezhad, T., Jamali, D.H., Hmida, A., Ozgoli, H.A., Safari,

- 587 S., AlShabi, M., Bani-Hani, E.H., 2022. A critical review of biogas production
- and usage with legislations framework across the globe. Int. J. Environ. Sci.
- 589 Technol. 19, 3377–3400. https://doi.org/10.1007/s13762-021-03301-6
- 590 Bautista Angeli, J.R., LeFloc'h, T., Lakel, A., Lacarrière, B., Andres, Y., 2022.
- 591 Anaerobic digestion of urban wastes: integration and benefits of a small-scale
- 592 system. Environ. Technol. 43, 3414–3425.
- 593 https://doi.org/10.1080/09593330.2021.1921857
- 594 Brockwell, P.J., Davis, R.A., 2016a. Modeling and Forecasting with ARMA Processes,
- in: Introduction to Time Series and Forecasting. Springer International
- ⁵⁹⁶ Publishing, Cham, pp. 121–155. https://doi.org/10.1007/978-3-319-29854-2_5
- 597 Brockwell, P.J., Davis, R.A., 2016b. Stationary Processes, in: Introduction to Time
- 598 Series and Forecasting. Springer International Publishing, Cham, pp. 39–71.

599 https://doi.org/10.1007/978-3-319-29854-2_2

- Brockwell, P.J., Davis, R.A., 2016c. ARMA Models, in: Introduction to Time Series
- and Forecasting. Springer International Publishing, Cham, pp. 73–96.
- 602 https://doi.org/10.1007/978-3-319-29854-2_3
- 603 Cowpertwait, P.S.P., Metcalfe, A. V, 2009. Time Series Data, in: Introductory Time
- 604 Series with R. Springer New York, New York, NY, pp. 1–25.
- 605 https://doi.org/10.1007/978-0-387-88698-5_1

606	James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. Statistical Learning, in: An
607	Introduction to Statistical Learning: With Applications in R. Springer New York,
608	New York, NY, pp. 15–57. https://doi.org/10.1007/978-1-4614-7138-7_2
609	James, G., Witten, D., Hastie, T., Tibshirani, R., Taylor, J., 2023. Multiple Testing, in:
610	An Introduction to Statistical Learning: With Applications in Python. Springer
611	International Publishing, Cham, pp. 557–596. https://doi.org/10.1007/978-3-
612	031-38747-0_13
613	Jeong, K., Abbas, A., Shin, J., Son, M., Kim, Y.M., Cho, K.H., 2021. Prediction of
614	biogas production in anaerobic co-digestion of organic wastes using deep
615	learning models. Water Res. 205, 117697.
616	https://doi.org/https://doi.org/10.1016/j.watres.2021.117697
617	Nguyen, L.N., Nguyen, A.Q., Nghiem, L.D., 2019. Microbial Community in
618	Anaerobic Digestion System: Progression in Microbial Ecology, in: Bui, XT.,
619	Chiemchaisri, C., Fujioka, T., Varjani, S. (Eds.), Water and Wastewater
620	Treatment Technologies. Springer Singapore, Singapore, pp. 331-355.
621	https://doi.org/10.1007/978-981-13-3259-3_15
622	Sadoune, H., Rihani, R., Marra, F.S., 2023. DNN model development of biogas
623	production from an anaerobic wastewater treatment plant using Bayesian
624	hyperparameter optimization. Chem. Eng. J. 471, 144671.
625	https://doi.org/https://doi.org/10.1016/j.cej.2023.144671
626	Sappl, J., Harders, M., Rauch, W., 2023. Machine learning for quantile regression of
627	biogas production rates in anaerobic digesters. Sci. Total Environ. 872, 161923.

628	https://	/doi.org/https:	//doi.org/	/10.1016/	′j.scitotenv	v.2023.161923
-----	----------	-----------------	------------	-----------	--------------	---------------

- 629 Shumway, R.H., Stoffer, D.S., 2017. ARIMA Models, in: Time Series Analysis and Its
- 630 Applications: With R Examples. Springer International Publishing, Cham, pp.
- 631 75–163. https://doi.org/10.1007/978-3-319-52452-8_3
- 632 Wang, Y., Huntington, T., Scown, C.D., 2021. Tree-Based Automated Machine
- 633 Learning to Predict Biogas Production for Anaerobic Co-digestion of Organic
- 634 Waste. ACS Sustain. Chem. Eng. 9, 12990–13000.
- https://doi.org/10.1021/acssuschemeng.1c04612
- 636 Yildirim, O., Ozkaya, B., 2023. Prediction of biogas production of industrial scale
- anaerobic digestion plant by machine learning algorithms. Chemosphere 335,
- 638 138976. https://doi.org/https://doi.org/10.1016/j.chemosphere.2023.138976