Natural Hazards
and Earth System
Sciences

# Meteorological Drought Trend Analysis and Forecasting Using a Hybrid SG-CEEMDAN-ARIMA-LSTM Model Based on SPI from Rain Gauge Data

**Siphamandla Sibiya**[1,3], **Shaun Ramroop**[1], **Sileshi Melesse**[1], and **Nkanyiso Mbatha**[2,4]

[1]School of Mathematics, Statistics, and Computer Science, University of KwaZulu-Natal, Pietermaritzburg Campus, Private Bag X01, Scottsville 3209, South Africa
[2]Council for Scientific and Industrial Research, Holistic Climate Change, Smart Places, Mering Naude Road, Brummeria, Pretoria 0001, South Africa
[3]Pure and Applied Analytics, Faculty of Natural and Agricultural Sciences, North-West University, Private Bag X1290, Potchefstroom, 2520, South Africa
[4]Department of Geography and Environmental Studies, University of Zululand, Private Bag X1001, KwaDlangezwa 3886, South Africa

**Correspondence:** Siphamandla Sibiya (siphasibiya@gmail.com)

**Abstract.** Meteorological drought presents considerable challenges to water supplies, agriculture, and socio-economic stability, especially in areas heavily reliant on precipitation. The Standardized Precipitation Index (SPI) is esteemed for its efficacy in drought monitoring, owing to its straightforwardness and applicability across many time scales. This study examines meteorological drought dynamics in the uMkhanyakude district using the Standardized Precipitation Index (SPI) at 6-, 9-, and 12-month timescales. Trend analysis was conducted using Mann–Kendall (MK), Modified Mann–Kendall (MMK), and Innovative Trend Analysis (ITA) methods. The study also proposes a hybrid model that integrates the Savitzky–Golay (SG) filter, Complete Ensemble Empirical Mode Decomposition with Adaptive Noise (CEEMDAN), Autoregressive Integrated Moving Average (ARIMA), and Long Short-Term Memory (LSTM) networks, referred to as SG-CEEMDAN-ARIMA-LSTM, for forecasting of the SPI time series. Analysis of SPI trends and variability revealed statistically significant declining trends at five monitoring stations, characterized by negative $Z$-scores and $p$-values, showing a marked downward trajectory across several SPI scales. On the other hand, the forecasting results demonstrate that the SG-CEEMDAN-ARIMA-LSTM methodology outper-formed benchmark models across all temporal scales, achieving high prediction accuracy with $R^2$ values of 0.9839 (SPI-6), 0.9892 (SPI-9), and 0.9990 (SPI-12). These findings highlight the effectiveness of decomposition techniques (SG, CEEMDAN) in enhancing model performance and confirm the suitability of the hybrid model for both short-term and long-term drought forecasting. This study merges robust trend analysis with advanced hybrid forecasting techniques, providing a reliable framework for early warning systems and sustainable water resource management in drought-prone regions.

## 1 Introduction

Drought is a complex and recurring natural hazard with significant economic, social, and environmental implications globally (Bagmar and Khudri, 2021; Kalisa et al., 2021; Song and Park, 2021). In contrast to other natural disasters, droughts manifest gradually, often persisting for extended periods, and their effects permeate various sectors, including agriculture, water resources, and socio-economic systems (Wilhite and Glantz, 1985; Cunha et al., 2019). This study specifically focuses on meteorological drought, char-

acterized as a sustained period of below-average precipitation (Taylan, 2024). Meteorological drought often serves as the initial phase that subsequently evolves into agricultural, hydrological, and socioeconomic drought (Malik et al., 2021; Latifoğlu and Özger, 2023). As it is solely influenced by precipitation variability, meteorological drought can be effectively quantified using precipitation-based indices.

Several indices have been established to quantify drought conditions, including the Standardized Precipitation Index (SPI) and the Standardized Precipitation Evapotranspiration Index (SPEI). While the SPEI integrates both precipitation and temperature data, its requirement for extensive datasets and complex computations may restrict its applicability in regions with limited data availability (Xu et al., 2020). Conversely, the SPI depends exclusively on precipitation, rendering it widely used for analysing meteorological drought, especially in semi-arid regions. Its versatility across multiple timescales facilitates the robust identification of both short- and long-term drought patterns. Accordingly, given the data constraints in the uMkhanyakude district of South Africa, this study adopts the SPI as the primary drought index, while recognizing that its exclusive reliance on precipitation constitutes a methodological limitation. Since SPI is precipitation-driven, analysing rainfall trends is a necessary first step before applying SPI under climate change conditions. Without first establishing rainfall trends, one risks misinterpreting SPI signals as short-term anomalies when they may actually reflect long-term climate-driven shifts.

In this context, the escalating concerns regarding climate change and its influence on local climates have underscored the necessity of analyzing drought trends. Thus, trend analysis of rainfall and SPI together provides a comprehensive picture of rainfall trends, revealing the climatic forcing, while SPI trends quantify the standardized drought intensity and persistence, which is crucial for understanding drought risk in the context of climate change. Systematic evaluations of drought occurrences not only contribute to the development of evidence-based water resource management strategies but also enhance the calibration of early warning systems and inform climate adaptation policies at both regional and national levels. Furthermore, temporal analyses enable researchers to assess the effectiveness of mitigation measures and anticipate emerging risks, thereby bolstering resilience in vulnerable sectors such as agriculture and public water supply. In the absence of structured trend analyses, drought management remains predominantly reactive, constraining the transition towards proactive and sustainable adaptation strategies. Building on trend analysis, drought forecasting is essential for deepening the understanding of drought dynamics. Effective forecasting provides early warnings that are critical for mitigating impacts and strengthening drought management strategies (Balti et al., 2020; Zhang et al., 2022, 2024; Tan et al., 2023).

Accurate forecasting of the SPI is crucial in regions such as uMkhanyakude, which is prone to recurrent and severe drought events. Enhanced prediction capabilities support agricultural resilience, water resource planning, and the establishment of early warning systems (Xu et al., 2020). Traditional statistical models, such as ARIMA or SARIMA, alongside contemporary machine learning methods, have been extensively employed for forecasting drought indices, including the SPI. However, each approach has inherent limitations. For example, Gudko et al. (2025) utilized SARIMA to analyze precipitation dynamics in Russia, demonstrating efficacy in short-term predictions while exhibiting constrained accuracy for long-term forecasts. Similarly, Hussain et al. (2025) integrated ARIMA with machine learning models to enhance SPI and SPEI predictions, achieving accuracies exceeding 92 %. This highlights the advantages of combining statistical and machine learning techniques. Nonetheless, these methodologies often encounter challenges associated with nonlinear and complex rainfall patterns, particularly over short time scales. To mitigate the limitations of standalone models, hybrid approaches have gained prevalence, capitalizing on the complementary strengths of diverse techniques. Alquraish et al. (2021) compared hybrid models such as HMM-GA, ARIMA-GA, and ARIMA-GA-ANN against, such as HMM-GA, ARIMA-GA, and ARIMA-GA-ANN, with conventional HMM and ARIMA models for SPI prediction in the Arabian Peninsula, revealing that hybrid models consistently outperformed their standalone counterparts. Likewise, Xu et al. (2022) and Ding et al. (2022) demonstrated that the combination of CEEMD with ARIMA or LSTM significantly improves SPI forecasts across multiple time scales in China, suggesting that decomposition-based hybrid methods effectively capture intricate temporal patterns.

Recent studies have significantly advanced hybrid methodologies through the implementation of sophisticated preprocessing and optimization techniques. Latifoğlu and Özger (2023) utilized phase transfer entropy (pTE) in conjunction with Tunable Q Factor Wavelet Transform (TQWT), optimized via Grey Wolf Optimization (GWO), followed by artificial neural networks (ANN), support vector regression (SVR), machine learning (ML), and Gaussian process regression (GPR), resulting in superior predictive performance. Sibiya et al. (2024) introduced the CEEMDAN-ARIMA-LSTM model for SPI predictions in Cape Town, demonstrating that the combination of CEEMDAN decomposition with both linear and nonlinear models can significantly improve forecast accuracy. Wei et al. (2025) adopted the Informer model and developed the VMD-JAYA-Informer hybrid, which integrates Variational Mode Decomposition (VMD) with an optimization algorithm, thereby enhancing short-term Standardized Precipitation Index (SPI) and Standardized Precipitation-Evapotranspiration Index (SPEI) forecasts.

Despite the successes achieved by hybrid models, several challenges persist. Decomposition techniques such as Empirical Mode Decomposition (EMD), Ensemble Em-

pirical Mode Decomposition (EEMD), Complete Ensemble Empirical Mode Decomposition with Adaptive Noise (CEEMDAN), and Variational Mode Decomposition (VMD) are computationally demanding, particularly when applied to large datasets or in real-time contexts (Sibiya et al., 2024). CEEMDAN, specifically, can yield misleading intrinsic mode functions (IMFs) when utilized on excessively noisy or unstable time series, which undermines the efficiency and reliability of subsequent predictions. Furthermore, existing research has not investigated the synergistic application of advanced smoothing filters in conjunction with decomposition techniques to mitigate noise prior to hybrid modeling.

To address these limitations, this study proposes an innovative hybrid model that integrates the Savitzky-Golay (SG) filter with CEEMDAN for preprocessing, followed by the Autoregressive Integrated Moving Average (ARIMA) and Long Short-Term Memory (LSTM) models for drought prediction. The SG filter is effective in smoothing high-frequency noise, thereby enhancing the decomposition process and alleviating the computational burden. The integration of the Savitzky-Golay smoothing filter with CEEMDAN substantially improves forecasting accuracy by enhancing the quality and interpretability of the input time series prior to modeling. This combination enables CEEMDAN to produce IMFs that are cleaner, more distinct, and less prone to spurious fluctuations, thus offering a more reliable foundation for subsequent predictive modeling. Cleaner IMFs facilitate the training of both linear (ARIMA) and nonlinear (LSTM) models, resulting in more accurate and robust forecasts. This approach capitalizes on the complementary strengths of both statistical and machine learning models, addressing noise-related issues inherent in raw data.

Although hybrid models have demonstrated superior performance in drought forecasting, no prior study has examined:

1. The combined use of smoothing techniques (SG filter) with CEEMDAN to enhance the quality of decomposition.

2. The implementation of an integrated SG-CEEMDAN-ARIMA-LSTM framework for trend-based Standardized Precipitation Index (SPI) predictions (SPI-6, SPI-9, SPI-12).

3. Forecasting efforts that explicitly incorporate both trend analysis and predictive modeling for semi-arid regions characterized by limited meteorological data.

As a result, the proposed SG-CEEMDAN-ARIMA-LSTM model addresses these gaps by enhancing decomposition efficiency, reducing computational costs, and improving prediction accuracy across multiple SPI timescales. This methodology offers valuable insights for water resource management, infrastructure planning, early warning systems, and the advancement of hybrid drought prediction models.

## 2 Material Methods

This study utilizes various time series forecasting models to analyse the intricate dynamics of meteorological drought as indicated by the Standardized Precipitation Index (SPI). The foundational statistical model examined is the Autoregressive Integrated Moving Average (ARIMA), which is adept at addressing linear relationships in time series data. The Long Short-Term Memory (LSTM) neural network is employed to tackle nonlinear patterns, supplemented by a hybrid ARIMA-LSTM framework that amalgamates the advantages of both models. Additional improvements are investigated by incorporating a Savitzky-Golay (SG) digital smoothing filter, which is often used to remove noise from time series or spectral data, into the ARIMA-LSTM model, and by utilizing the Complete Ensemble Empirical Mode Decomposition with Adaptive Noise (CEEMDAN) before ARIMA-LSTM to more effectively manage nonstationary signals. The work introduces a unique hybrid model, SG-CEEMDAN-ARIMA-LSTM, which integrates decomposition and hybrid modeling techniques to enhance the accuracy and robustness of drought forecasts.

Therefore, the subsequent Materials and Methods section will provide a detailed account of the study area, the data employed, and the preprocessing steps undertaken, including the trend extraction methods applied prior to forecasting. This will be followed by an in-depth description of each modeling approach, outlining their theoretical foundations, implementation procedures, and parameterization strategies. Such a structured presentation ensures transparency in model development and establishes a comprehensive methodological framework for the proposed forecasting system.

### 2.1 Study Area and Data

This study employed monthly mean precipitation records from 1980 to 2023, obtained from the South African Weather Service (SAWS) for the uMkhanyakude District in South Africa. The uMkhanyakude District Municipality is located in the far northern region of the KwaZulu-Natal (KZN) province (coordinates: 32.014489° S, 27.622242° E). The municipality covers a total area of 13 855 km$^2$, making it the second largest in the province, exceeded only by the Zululand Municipality. The uMkhanyakude District was formed immediately after the local government elections in December 2000, as part of the municipal demarcation process, encompassing some of the most destitute and underdeveloped areas of KwaZulu-Natal. The uMkhanyakude District consists of four local municipalities: uMhlabuyalingana, Jozini, Big Five Hlabisa, and Mtubatuba. The municipality is geographically surrounded by Mozambique to the north, the Indian Ocean to the east, the uThungulu River to the south, Zululand to the west, and the Kingdom of Swaziland to the northwest. Figure 1 illustrates the spatial distribution of the stations.
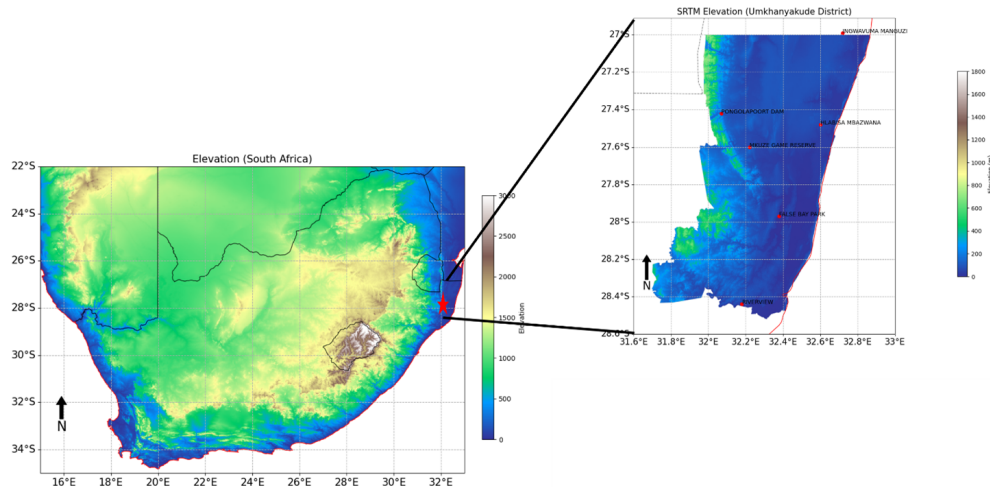
**Figure 1.** Overview of the uMkhanyakude District, South Africa. Rain gauge stations are marked red.
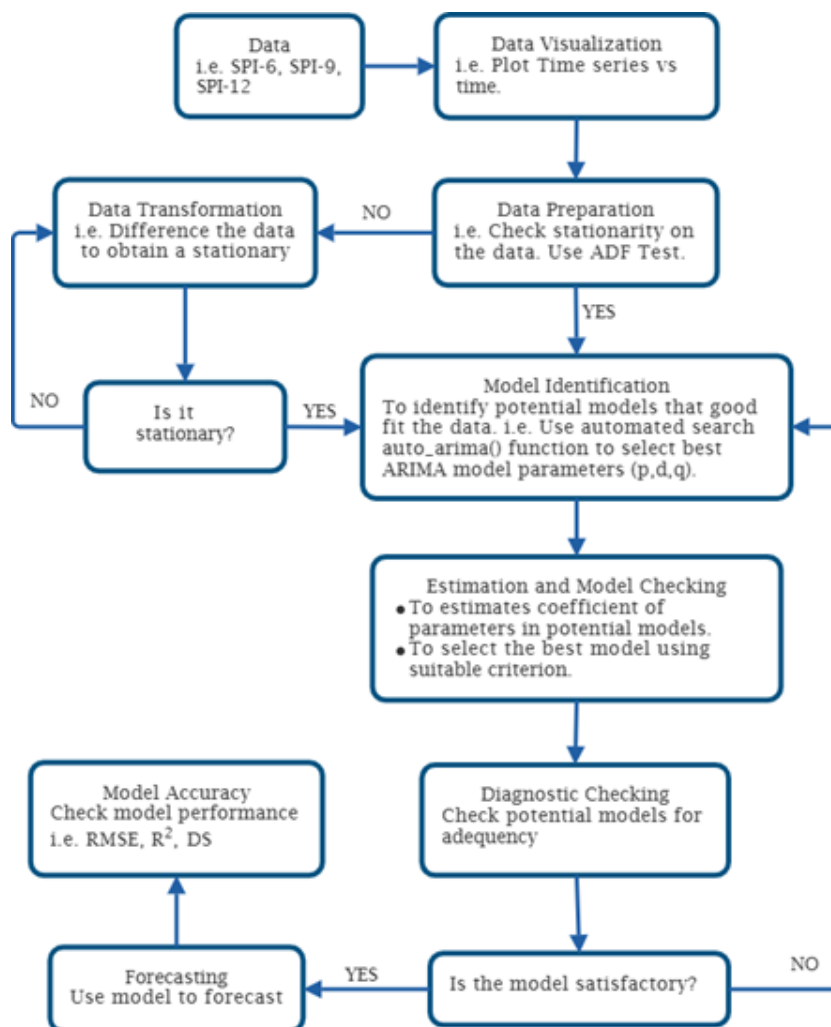


**Figure 2.** The Box-Jenkins Steps Approach.

## 2.2 Modified Mann-Kendall

The modified Mann-Kendall methodology is derived from the nonparametric Mann-Kendall method (Mann, 1945; Kendall, 1975), which is widely used to detect trends in hydro-meteorological time series (Caloiero et al., 2011; Bard et al., 2015; Wang et al., 2017; Mirabbasi et al., 2020). The modified Mann–Kendall (MMK) test was employed for serially correlated data exhibiting a substantial lag-1 autocorrelation coefficient, utilising the variance correction method proposed by Yue et al. (2002). Hamed and Rao (1998) created this methodology to eradicate all substantial autocorrelation in the time series. Under the assumption that the data are independent and identically distributed, the $S$ statistic of the Mann-Kendall test is computed as follows (Sharifi et al., 2024):

$$S = \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} \mathrm{Sign}\left(x_j - x_i\right) \tag{1}$$

where $n$ denotes the sample size; $x_i$ and $x_j$ denote sequential $i$th and $j$th data points, respectively, and sign(.) is the sign function which can be computed as

$$\mathrm{Sign}\left(x_j - x_i\right) = \begin{cases} 1, & \text{if } x_j - x_i > 0 \\ 0, & \text{if } x_j - x_i = 0 \\ -1, & \text{if } x_j - x_i < 0 \end{cases} \tag{2}$$

with the mean and variance of the $S$ statistics in the equation are as follows (Helsel and Hirsch, 1993; Ma et al., 2014; Ashraf et al., 2023)

$$E(S) = 0 \tag{3}$$

$$\mathrm{Var}(S) = \frac{n(n-1)(2n+5) - \sum_{i=1}^{p} t_i(t_i-1)(2t_i+5)}{18} \tag{4}$$

where $p$ is the number of tied groups and $t_i$ denotes the number of data points in the $t$th group. The second term represents an adjustment for tied group or censored data. The standardized Z statistic is calculated as

$$Z_{\mathrm{MK}} = \begin{cases} \frac{S-1}{\sqrt{\mathrm{Var}(S)}}, & S > 0 \\ 0, & S = 0 \\ \frac{S+1}{\sqrt{\mathrm{Var}(S)}}, & S < 0 \end{cases} \tag{5}$$

The test statistic $Z$ is used to measure the significance of the trends. In the modified Mann-Kendall approach, a modified variance of $S$ is computed as follows (Hamed and Rao, 1998)

$$\mathrm{Var}(S^*) = \mathrm{Var}(S) \frac{n}{n^*} \tag{6}$$

where $n^*$ is the effective sample size. The $\frac{n}{n^*}$ ratio can be calculated as follows (Hamed and Rao, 1998)

$$\frac{n}{n^*} = 1 + \frac{2}{n(n-1)(n-2)} \sum_{i=1}^{n} (n-i)$$
$$(n-i-1)(n-i-2)\, r_i \tag{7}$$

where $r_i$ denotes the lag-$i$ significant autocorrelation coefficient of rank $i$ in a time series. Then the standardized statistic of the S statistic, denoted as $Z$, can be derived as

$$Z_{\mathrm{MMK}} = \begin{cases} \frac{S-1}{\sqrt{\mathrm{Var}(S^*)}}, & S > 0 \\ 0, & S = 0 \\ \frac{S+1}{\sqrt{\mathrm{Var}(S^*)}}, & S < 0 \end{cases} \tag{8}$$

If the calculated $Z$ values ($Z_{\mathrm{MK}}$ and $Z_{\mathrm{MMK}}$) exceed the critical values of $-Z_{1-\alpha/2}$ or fall below $Z_{1-\alpha/2}$, there is no discernible trend in the time series at the significance level of $\alpha$. If the $Z$ value is positive and exceeds $Z_{1-\alpha/2}$, the trend is upward; conversely, if the $Z$ value is negative and falls below $-Z_{1-\alpha/2}$, the trend is downward.

## 2.3 Innovative Trend Analysis

The Innovative Trend Analysis (ITA) method, initially introduced by Şen (2012), has been widely employed for detecting patterns in precipitation time series. Since its debut, the ITA technique has experienced substantial improvements in both mathematical and graphical aspects, as evidenced by Şen (2017) and Alashan (2018). The ITA method does not depend on assumptions of serial autocorrelation, normalcy, or record length, making it appropriate for both graphical and statistical trend analysis (Zena et al., 2022). Initially, the time series is bifurcated into two equal segments and organised in ascending order. The initial segment of the time series ($x_i : i = 1, 2, \ldots, n/2$) is positioned along the horizontal $x$-axis, while the subsequent segment ($x_j : j = n/2+1, n/2+2, \ldots, n$) is situated along the vertical $y$-axis in the Cartesian coordinate system (Ashraf et al., 2023). The ITA approach visually represents trend analysis, specifically indicating monotonic growing, declining, and trendless circumstances (Öztopal and Şen, 2017; Likinaw et al., 2023). A monotonically growing or declining trend can be identified when the majority of points are situated above or below the 45° (1 : 1 line), respectively. A trendless condition arises when the data points are clustered along the 45° line (Şen, 2012). We employ the magnitude of the slope parameter to convey information about monotonicity. The slope parameter of the ITA technique is a stochastic property dependent on the sample means of the first half ($n_1$) and the second half ($n_2$) of the time-series mean data values. According to Şen (2017), the straight-line trend slope ($S_{\mathrm{ITA}}$) can be estimated using the following expression:

$$S_{\mathrm{ITA}} = \frac{2x\left(x_j - x_i\right)}{n} \tag{9}$$

where $n$ represents the total number of observations, $x_i$ and $x_j$ are the arithmetic means of the first and second halves of the sub-series, respectively. Given that $x_i$ and $x_j$ are stochastic variables, the expected value of the slope can be determined by analysing the expectancies of both the first and second halves of the time series (Alashan, 2020; Harka et
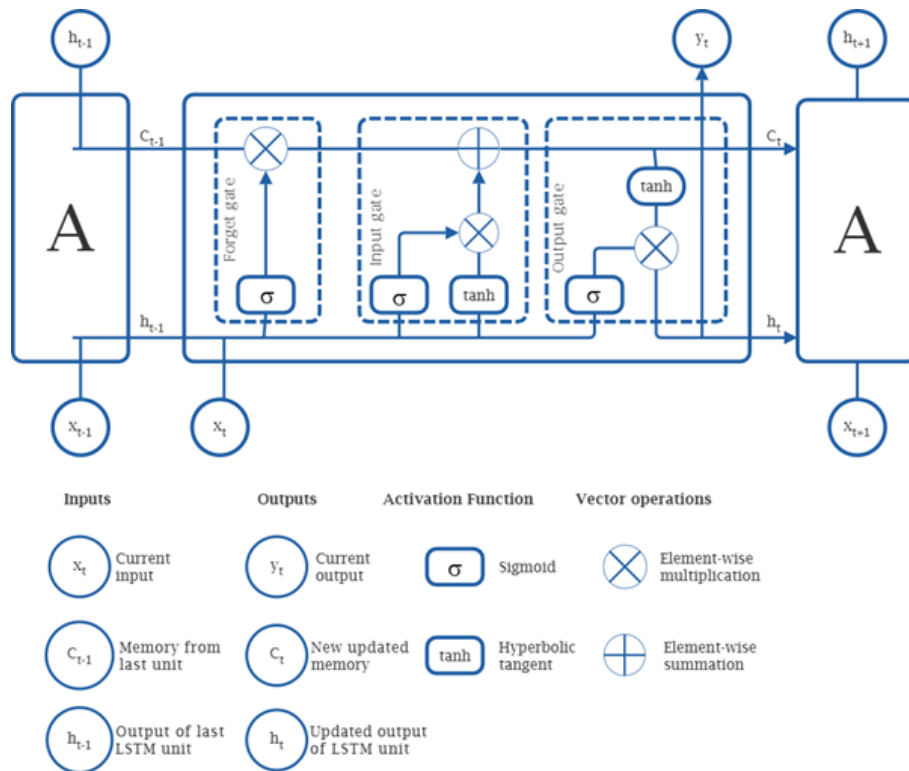
**Figure 3.** Structure diagram of LSTM model.

al., 2021):

$$E(S_{\text{ITA}}) = \frac{2}{n}\left[E(x_j) - E(x_i)\right] \qquad (10)$$

For the no trend condition, $E(x_j) = E(x_i)$, the $E(S_{\text{ITA}}) = 0$ and standard deviation (SD) of the two half time-series $(\sigma_{x_j} = \sigma_{x_i} = \sigma/\sqrt{n})$, $\sigma$ is the SD is of the parent series. If $E(x_j) \neq E(x_i)$, the differences between $E(x_j)$ and $E(x_i)$ gives the variance

$$\sigma^2_{S_{\text{ITA}}} = \frac{8}{n^2}\left[E(x_j) - E(x_j x_i)\right] \qquad (11)$$

and the SD of the slope

$$\sigma_{S_{\text{ITA}}} = \frac{2\sqrt{2}}{n\sqrt{n}}\sigma\sqrt{1 - \rho_{x_j x_i}} \qquad (12)$$

In the stochastic processes, the term $\rho_{x_j x_i}$ is the correlation coefficient between the two mean values, and can be estimated as

$$\rho_{x_j x_i} = \frac{E(x_j x_i) - E(x_j)E(x_i)}{\sigma_{x_j}\sigma_{x_i}} \qquad (13)$$

In the end, the upper and lower confidence limit (CL) of the trend slope was calculated (Şen, 2017):

$$\text{CL}_{(1-\alpha)} = 0 \pm (Z_{1-\alpha/2})\sigma_{S_{\text{ITA}}} \qquad (14)$$

$Z_{1-\alpha/2}$ denotes the crucial slope for standardised time-series at $\pm 1.96$ for a 95 % significance level or $\pm 1.645$ for a 90 % significance level (Alashan, 2020). If the ITA slope value is beyond the lower and upper confidence limits, the null hypothesis of no significant trend should be rejected at the $\alpha$ significance level (Şen, 2017). In a two-tailed scenario, the null hypothesis ($H_0$) posits the absence of a trend in time-series data, while the alternative hypothesis ($H_1$) asserts the presence of a trend in time-series data at a significance level of $\alpha$. If the slope, $\pm S_{\text{ITA}} > \pm\text{CL}_{(1-\alpha)}$, then ($H_0$) is discarded in favour of ($H_1$). The positive and negative values of $S_{\text{ITA}}$ signify an upward and downward trend in the time-series data, respectively (Şen, 2017).

### 2.4 The SPI Calculation

For the purpose of analysing the severity of drought, which is caused by a lack of water supply as a result of reduced precipitation in response to rising demand, the SPI was created by McKee et al. (1993) and is based on probability (Zuo et al., 2021). Based on the cumulative likelihood of a specific amount of precipitation, the SPI indicator is calculated by fitting the precipitation throughout the same period with a certain distribution function. At its largest point, the SPI index represents the quantile of a normal distribution. Each time axis has an estimated drought index for 6, 9, and 12 months. This is based on the gamma probability density func-

tion, which accounts for the periodic distribution of precipitation for the corresponding data point. The expression of the density function for this distribution is as follows.

$$g(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\beta}} \quad (15)$$

where $\alpha$ is the shape parameter, $\beta$ is the scale parameter and $x$ is the precipitation amount, and $\Gamma(\alpha) = \int_0^\infty y^{\alpha-1} e^{-y} dy$ is gamma function. The maximum likelihood estimates of the parameters $\alpha$ and $\beta$ are:

$$\alpha = \frac{1}{4A}\left(1 + \sqrt{1 + \frac{4A}{3}}\right) \quad (16)$$

$$\beta = \frac{\overline{x}}{n} \quad (17)$$

where $A = \ln(\overline{x}) - \frac{\sum \ln(x)}{n}$, $\overline{x}$ is the precipitation average and $n$ is the sample size. The following equation applies the acquired parameters to the cumulative probability distribution:

$$G(x) = \int_0^x g(x)\, dx = \frac{1}{\beta^\alpha \Gamma(\alpha)} \int_0^x x^{\alpha-1} e^{-\frac{x}{\beta}} dx \quad (18)$$

$G(x)$ represents the likelihood that the precipitation will be equal to or less than $x$. The distribution function for precipitation needs to be adjusted because the real precipitation samples can contain a value of 0. Based on this, we can calculate the cumulative probability as:

$$H(x) = q + (1-q) G(x) \quad (19)$$

where $q$ denotes the probability when precipitation equals zero. The probability of no rainfall, $q$, can be articulated as $q = m/r$, where m represents the number of days without rainfall and r denotes the number of days with rainfall (Song and Park, 2021). Consequently, $H(x)$ is converted to the conventional random variable $Z$ of the standard normal distribution, characterised by a mean of 0 and a variance of 1, resulting in:

$$\mathrm{SPI} = Z = \begin{cases} -\left(k - \frac{c_0 + c_1 k + c_2 k^2}{1 + d_1 k + d_2 k^2 + d_3 k^3}\right), & 0 < H(x) \leq 0.5 \\ +\left(k - \frac{c_0 + c_1 k + c_2 k^2}{1 + d_1 k + d_2 k^2 + d_3 k^3}\right), & 0 < H(x) \leq 1.0 \end{cases} \quad (20)$$

$$k = \begin{cases} \sqrt{\ln\left(\left(\frac{1}{H(x)}\right)^2\right)}, & 0 < H(x) < 0.5 \\ \sqrt{\ln\left(\left(\frac{1}{1-H(x)}\right)^2\right)}, & 0 < H(x), < 1.0 \end{cases} \quad (21)$$

where $c_0 = 2,515517$, $c_1 = 0.802853$, $c_2 = 0,010328$, $d_1 = 1,432788$, $d_2 = 0,189269$, $d_3 = 0,001308$ are constants. Furthermore, the SPI indicator is a standardised normalised index, establishing a correlational relationship with likelihood. Table 1 presents the probability associated with each category of drought.

**Table 1.** Drought classification using SPI values and corresponding event probability (Lloyd-Hughes and Saunders, 2002).

| SPI Values | Drought Category | Probability (%) |
|---|---|---|
| $2.00 \leq \mathrm{SPI}$ | Extremely wet | 2.3 |
| $1.50 \leq \mathrm{SPI} \leq 1.99$ | Severely wet | 4.4 |
| $1.00 \leq \mathrm{SPI} \leq 1.49$ | Moderately wet | 9.2 |
| $0.00 \leq \mathrm{SPI} \leq 0.99$ | Mildly wet | 34.1 |
| $-0.99 \leq \mathrm{SPI} \leq 0.00$ | Mild dry | 34.1 |
| $-1.49 \leq \mathrm{SPI} \leq -1.00$ | Moderate dry | 9.2 |
| $-1.99 \leq \mathrm{SPI} \leq -1.50$ | Severe dry | 4.4 |
| $\mathrm{SPI} \leq -2.00$ | Extreme dry | 2.3 |

## 2.5 The Savitzky-Golay Filter

The Savitzky-Golay (SG) smoothing technique is a widely used method for noise filtration. Savitzky and Golay (1964) introduced the SG filter as an effective technique for signal smoothing. The SG technique attenuates noise utilising two parameters: polynomial order and window size. By flexibly adjusting these two parameters, the SG filter can achieve exceptional performance in various pre-processing circumstances. The essence of this procedure involves fitting a low-degree polynomial to the samples within a sliding window using the least squares method, resulting in a new smoothed value for the central point derived by convolution. The SG filter is a specific variant of a low-pass filter that substitutes each value in the time series with a new value derived from a polynomial fit to $2m + 1$ surrounding points, including the point to be smoothed, where m is equal to or larger than the polynomial's order. The polynomial is articulated as follows:

$$p(n) = \sum_{k=0}^{N} a_k n^k \quad (22)$$

where $N$ is the power of the polynomial and $N \leq 2M + 1$. The following equation is used to determine the error between the estimated and original values; in order to find the desired polynomial result, this error must be minimised.

$$\epsilon_N = \sum_{n=-M}^{M} (p(n) - x[n])^2 \quad (23)$$

The following form of discrete convolution can be used to express the filter's output:

$$y[n] = \sum_{m=-M}^{M} h[m]\, x[n-m] = \sum_{m=n-M}^{n+M} h[n-m]\, x[m] \quad (24)$$

This work employs the SG filter for two primary reasons: firstly, it enhances system performance by preserving the width and height of waveform peaks in noisy SPI, and secondly, it modifies the SPI while maintaining its fundamental qualities.
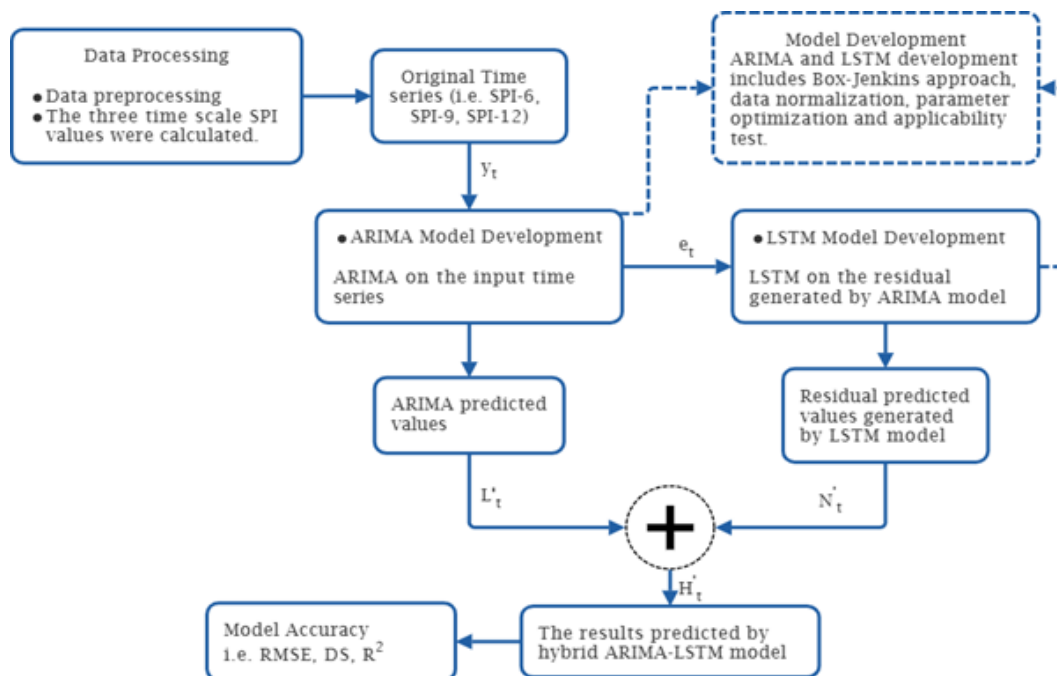
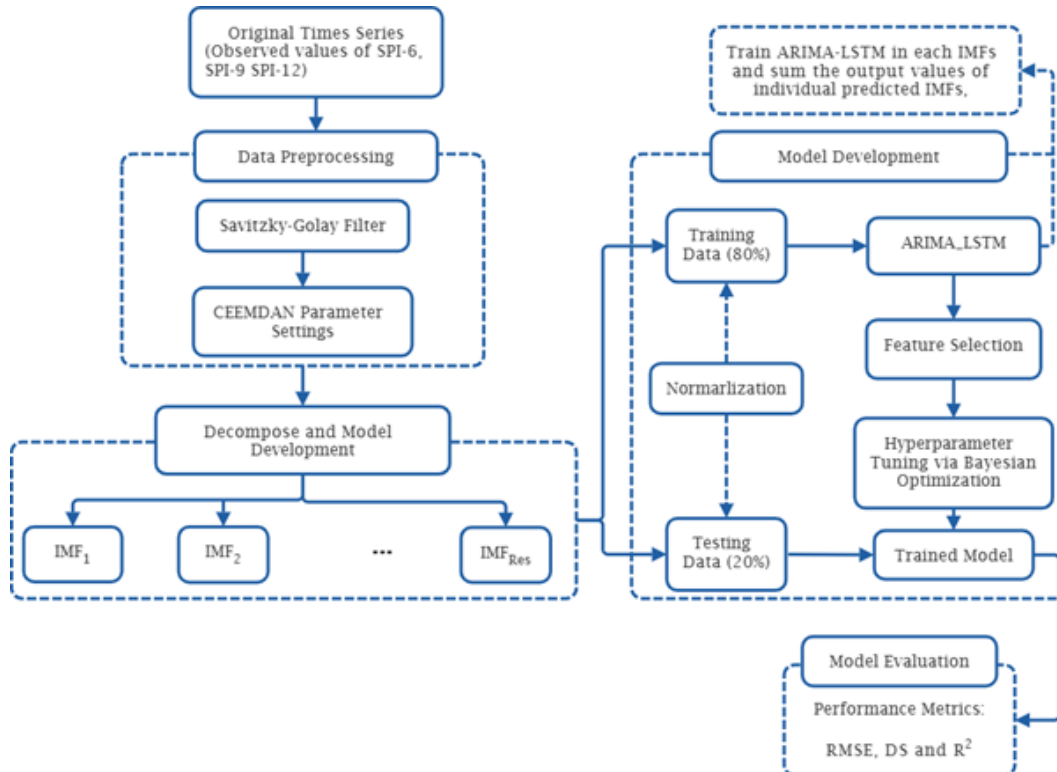**Figure 4.** Predictive flowchart of the ARIMA-LSTM hybrid model.

**Figure 5.** Procedure of proposed SG-CEEMDAN-ARIMA-LSTM hybrid model.

## 2.6 The Complete Ensemble Empirical Mode Decomposition with Adaptive Noise

The model's ability to fit functions and converge will be constrained by the complexity and volatility of the original time sequence, which in turn limits the model's predictive power. To overcome this challenge, the complete ensemble empirical mode decomposition (CEEMDAN) technique is employed to preprocess the original nonstationary and nonlinear time series. Both empirical mode decomposition (EMD) and ensemble empirical mode decomposition (EEMD), have been enhanced by the CEEMDAN. The computational efficiency is improved, and the reconstructed sequences of both the EMD and EEMD algorithms are free of modal confusion and noise residuals (Zhang et al., 2023). A residual term and a sequence of intrinsic mode functions (IMFs) are the building blocks of a complicated time series signal that the CEEMDAN breaks down.

Step 1: Incorporate a constrained quantity of adaptive white noise into the original sequence $x(t)\delta_0\omega^i(t)(t = 1, 2, 3, \ldots, N)$

$$x^i(t) = x(t) + \delta_0\omega^i(t) \tag{25}$$

where $N$ denotes the number of trials, $\delta_0$ signifies a coefficient of intensity, and $\omega^i(t)$ indicates the ith realisation of a stochastic Gaussian process.

Step 2: The residual $r_1(t)$ and the first modal component $\text{IMF}_1$ are obtained by decomposing each Eq. (1) using EMD.

$$\overline{\text{IMF}_1(t)} = \frac{1}{N}\sum_{i=1}^{N}\text{EMD}_1\left(x^i(t)\right) \tag{26}$$

$$r_1(t) = x(t) - \overline{\text{IMF}_1(t)} \tag{27}$$

In this context, $\text{EMD}_1(.)$ denotes the initial IMF component produced by the EMD algorithm, while $r_1(t)$ signifies the residual associated with the first stage.

Step 3: Add white noise $\delta_1\text{EMD}_1(\omega^i(t))$ to the residual $r_1(t)$ and further decomposed by EMD to obtain the second modal component $\text{IMF}_2$ and residual $r_2(t)$.

$$\overline{\text{IMF}_2(t)} = \frac{1}{N}\sum_{i=1}^{N}\text{EMD}_1\left(r_1(t) + \delta_1\text{EMD}_1(\omega^i(t))\right) \tag{28}$$

$$r_2(t) = r_1(t) - \overline{\text{IMF}_2(t)} \tag{29}$$

For the $j = 3, 4, \ldots, N$, the $j$th IMF component and the $j$th residual can be computed as:

$$\overline{\text{IMF}_j(t)} = \frac{1}{N}\sum_{i=1}^{N}\text{EMD}_1\left(r_{j-1}(t) + \delta_{j-1}\text{EMD}_{j-1}(\omega^i(t))\right) \tag{30}$$

$$r_j(t) = r_{j-1}(t) - \overline{\text{IMF}_j(t)} \tag{31}$$

where $\text{EMD}_{j-1}(.)$ denotes the $(j-1)$th intrinsic mode function component derived from the empirical mode decomposition technique, and $r_j(t)$ represents the residual following the jth decomposition.

Step 3: Continue executing step 3 until the residual $r_j(t)$ meets a predetermined termination criterion.

The time series $x(t)$ can ultimately be articulated as

$$x(t) = \sum_{i=1}^{N}\overline{\text{IMF}_N(t)} + r_N(t) \tag{32}$$

## 2.7 The Autoregressive Integrated Moving Average Model

The Autoregressive Integrated Moving Average (ARIMA) model, pioneered by Box and Jenkins in the 1970s, serves as a robust and effective forecasting approach for time series analysis (Box et al., 2015). The ARIMA model, often known as the Box-Jenkins approach, is depicted through the concepts presented by Sibiya et al. (2024) in Fig. 2. The ARIMA models predict future values of the time series as a linear combination of historical and residual data. This model comprises three components: the order of seasonal differentiation, autoregressive order, and moving average order (Montgomery et al., 2015). The backward shift operator $B$ is employed to eliminate nonstationarity. A time series, $y_t$, is called homogeneous nonstationary if it first order difference, $\omega_y = (1 - B)y_t = y_t - y_{t-1}$ or the $d$th difference $\omega_t = (1 - B)^d y_t$ is also stationary time series. Furthermore, $y_t$ is referred to as an ARIMA model with orders $pd$ and $q$, noted ARIMA($pdq$). Hence, an ARIMA($pdq$) is often expressed as

$$\phi(B)(1 - B)^d y_t = c + \theta(B)\varepsilon_t \tag{33}$$

$$\phi(B) = 1 - \sum_{i=1}^{p}\phi_i B^i \quad \text{and} \quad \theta(B) = 1 - \sum_{i=1}^{q}\theta_i B^i \tag{34}$$

The backward shift operators for $\text{AR}(p)$ and $\text{MA}(q)$ are defined as $\phi(B)y_t = c + \varepsilon_t$ and $y_t = \mu + \theta(B)\varepsilon_t$ with $c = \mu - \phi\mu$, where $\mu$ and $\varepsilon_t$ are the mean and white noise, respectively and the $\varepsilon_t$ is independent and normal distributed with mean and variance of $\sigma_\varepsilon^2$.

## 2.8 The Long Short-Term Memory

Long short-term memory (LSTM) algorithms represent a category of recurrent neural network (RNN) designs that are proficient in handling sequential input and identifying temporal relationships (Hochreiter and Schmidhuber, 1997). LSTM networks incorporate specific memory cells and gates for the efficient management and regulation of information flow over various time steps. Consequently, they can effectively represent the data input while maintaining essential dependencies and patterns. The LSTM methodology addresses the problem of vanishing gradients encountered by RNN algorithms. This occurs when the gradient diminishes to a level insufficient for effectively updating the weights throughout prolonged sequences. The LSTM facilitates the flow of gradients across time by employing memory cells and gates. The
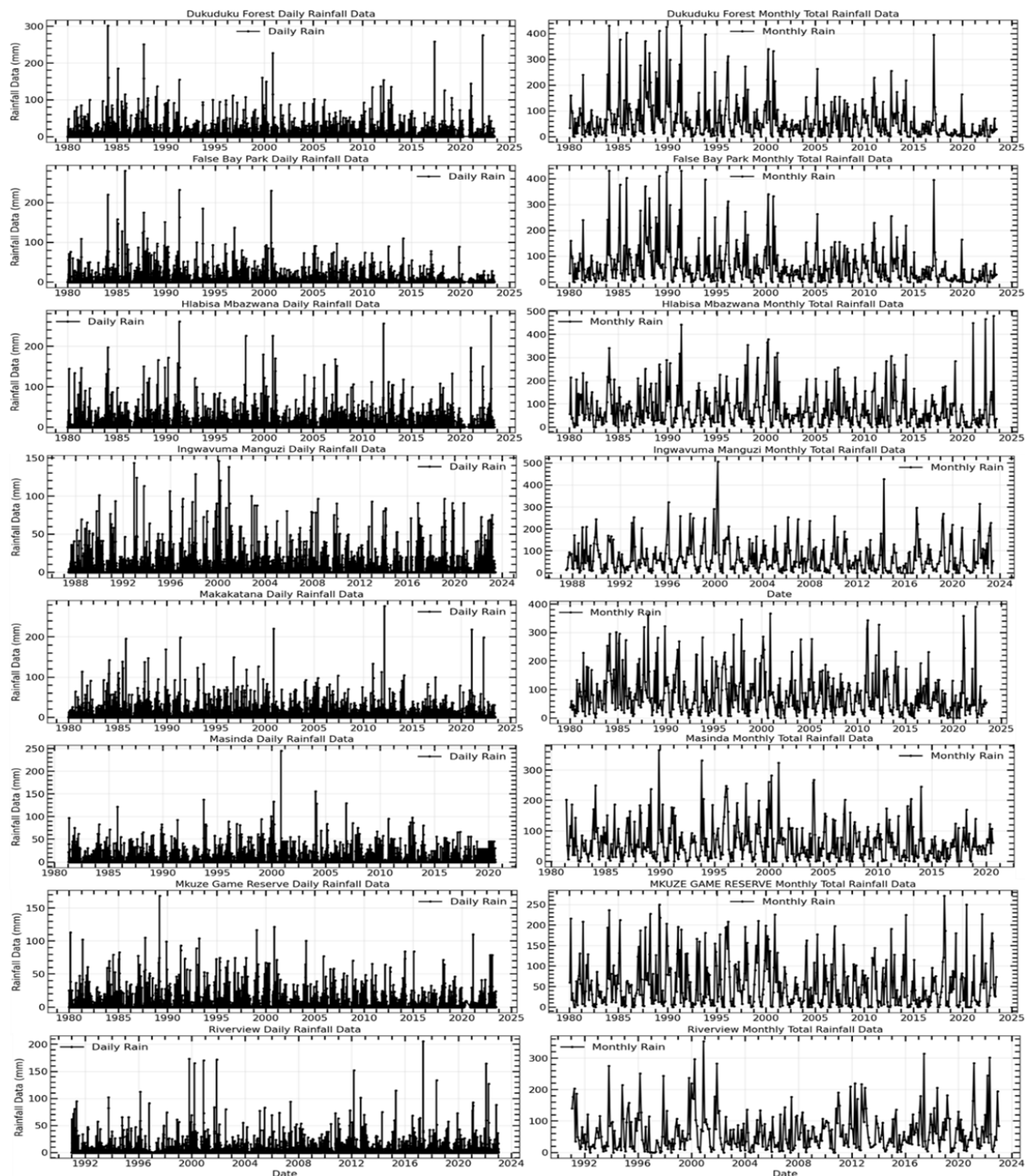
https://doi.org/10.5194/nhess-26-315-2026

Nat. Hazards Earth Syst. Sci., 26, 315–342, 2026

**Figure 6.** Time series plots of daily and monthly total rainfall data for uMkhanyakude district from early 1980's to 2023. The (left) plot shows the daily rainfall data in millimeters (mm), illustrating the high variability and intermittent nature of daily rainfall events over the years. The (right) plot presents the monthly total rainfall data (mm), which smooths out the daily variability and reveals clearer patterns of rainfall distribution over time.
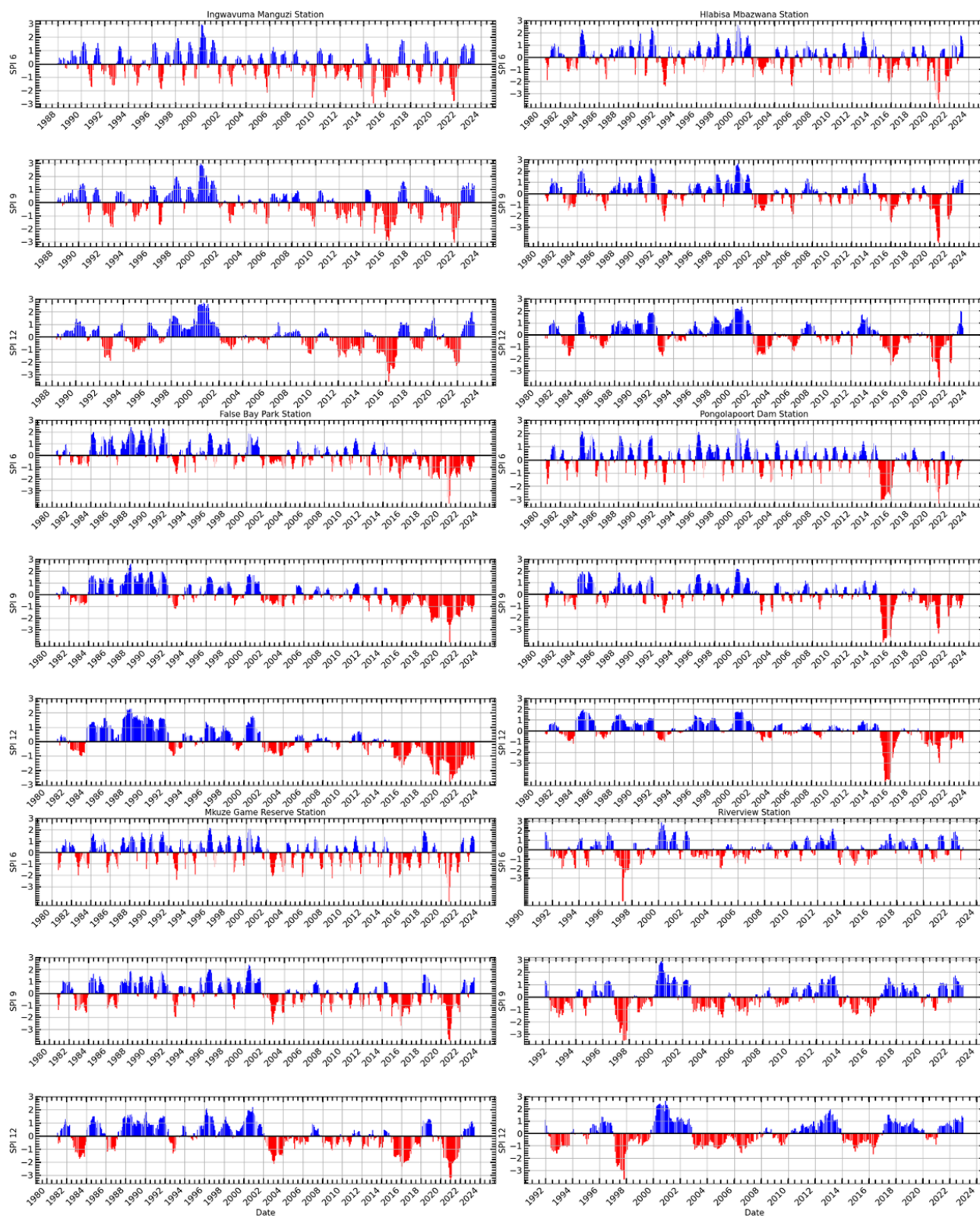
**Figure 7.** Standardized Precipitation Index (SPI) time series plots for uMkhanyakude district over 6-month (SPI-6), 9-month (SPI-9), and 12-month (SPI-12) periods from early 1980's to 2023. Positive SPI values (blue bars) indicate wetter-than-normal conditions, while negative SPI values (red bars) indicate drier-than-normal conditions.

model's foundational design primarily consists of three control gates: input, forget, and output. The activation function is represented by $\sigma$, whereas the cell states at time $t-1$ and $t$ are designated as $C_{t-1}$ and $C_t$ respectively. At time $t$ and time $t-1$, the cell possesses two concealed states, $h_t$ and $h_{t-1}$. Figure 3 illustrates the building of the LSTM unit, and the mathematical Eqs. (35) to (40) for the LSTM method are provided below. Initially, by employing the model's forget gate, we may determine the current hidden state $h_{t-1}$ and the degree to which the input $x_t$ has been preserved. The formula is

$$f_t = \sigma\left(W_f x_t + U_f h_{t-1} + b_f\right) \tag{35}$$

Secondly, the input gate allows us to ascertain the volume of content from the input variable that can be retained in the cell state $C_t$

$$i_t = \sigma\left(W_i x_t + U_i h_{t-1} + b_i\right) \tag{36}$$

$$\tilde{C}_t = \sigma_c\left(W_c x_t + U_i h_{t-1} + b_i\right) \tag{37}$$

$$C_t = f_t \odot C_{t-1} + i_i \odot \tilde{C}_t \tag{38}$$

The output gate of the LSTM produces outputs, and the hidden state of each cell is represented by the formula:

$$o_t = \sigma\left(W_o x_t + U_o h_{t-1} + b_o\right) \tag{39}$$

$$h_t = O_t \odot \sigma_h\left(C_t\right) \tag{40}$$

In the aforementioned formulas, $W_f$, $W_i$, and $W_o$ represent the weight matrices associated with the various control gates. The terms $b_f$, $b_i$, and $b_o$ correspond to the bias terms for each respective control gate. The notation $\tilde{C}_t$ signifies the complete input activation vector, while the operator $\odot$ (Hadamard product) indicates the element-wise multiplication of the elements between two vectors. The $\sigma$ activation function quantifies the amount of information that is transmitted through the various control gates.

## 2.9 The ARIMA-LSTM hybrid Model

Achieving accurate estimates of SPI index values through a forecasting model is essential for informed decision-making. Zhang (2003) offers a hybrid model wherein the ARIMA model extracts and predicts linear components, while the residuals, representing nonlinear data subcomponents, are then modelled by the LSTM approach. This study employs a hybrid model that integrates ARIMA and LSTM to predict both linear and nonlinear behaviours with optimal accuracy.

$$H_t = L_t + \aleph_t \tag{41}$$

where $L_t$ and $\aleph_t$ denote the linear and nonlinear components, respectively, for the hybrid technique which are computed using the initial time series ($y_t$). Consider the original dataset at time t and the forecast results obtained from applying the ARIMA model as $\hat{L}_t$ the prediction results. Thus,

$E_t = y_t - \hat{L}_t$ is the definition of the residual $E_t$ that is derived by removing $\hat{L}_t$ from $y_t$. Subsequently we compute the value $\hat{\aleph}_t$ by feeding the series of residuals into the LSTM model, which predicts the nonlinear component of the values. This equation may be written as

$$\hat{\aleph}_t = f_{\text{LSTM}}\left(E_{-1}, E_{-2}, \ldots, E_{-n}\right) + \epsilon_t, \tag{42}$$

where $\hat{\aleph}_t$ is a nonlinear expression associated with the LSTM model and $\epsilon_t$ is the random error. The combined forecasts from the two steps were then used to determine the value for the ARIMA-LSTM hybrid model. As illustrated in Fig. 4, the equation $\hat{H}_t = \hat{L}_t + \hat{\aleph}_t$ predicts the linearity and nonlinearity values, respectively, using ARIMA and LSTM models.

## 2.10 The development of the proposed SG-CEEMDAN-ARIMA-LSTM hybrid model

Due to the great uncertainty of the drought data and the existence of complexity, nonlinearity, and nonstationary trends, the single prediction model is greatly limited; however, the hybrid method has better prediction accuracy. The SG-CEEMDAN-ARIMA-LSTM algorithm that combines different techniques for improved accuracy in predicting drought based on the standardised precipitation index is proposed this study. This hybrid model is designed as a sequential framework where each step refines the data for subsequent modelling. The SG-CEEMDAN pre-processing stage enhances the data by smoothing and decomposing it into the meaningful components. The benefits of integrating the Savitzky–Golay smoothing filter with CEEMDAN significantly contribute to the enhancement of forecasting accuracy by improving the quality and interpretability of the input time series prior to modeling. The Savitzky–Golay filter acts as a noise suppression mechanism that preserves essential features of the time series, while eliminating high-frequency noise. This step ensures that the input to the CEEMDAN decomposition process is already denoised, leading to more stable and physically meaningful decomposed components. The CEEMDAN generates IMFs that are cleaner, more distinct, and less affected by spurious fluctuations. This results in better mode separation, reduces signal leakage across IMFs, and enhances the stationarity and regularity of each component. This hybrid preprocessing pipeline can enhances model generalization, reduces overfitting, and ultimately leads to more reliable and accurate forecasts. The components fed to the ARIMA-LSTM model that involves two-step process: the ARIMA for initial prediction utilising the Box-Jekins methodology and the LSTM model for refining and enhancing predictions. The hybrid model combines the ARIMA and the LSTM predictions to form the final hybrid forecasts. Figure 5 illustrates the proposed hybrid model algorithm. The process of SPI prediction based on ARIMA-LSTM combined with SG and CEEMDAN as is shown in Fig. 5. The process of the data smoothing, decomposition and prediction include four main steps.
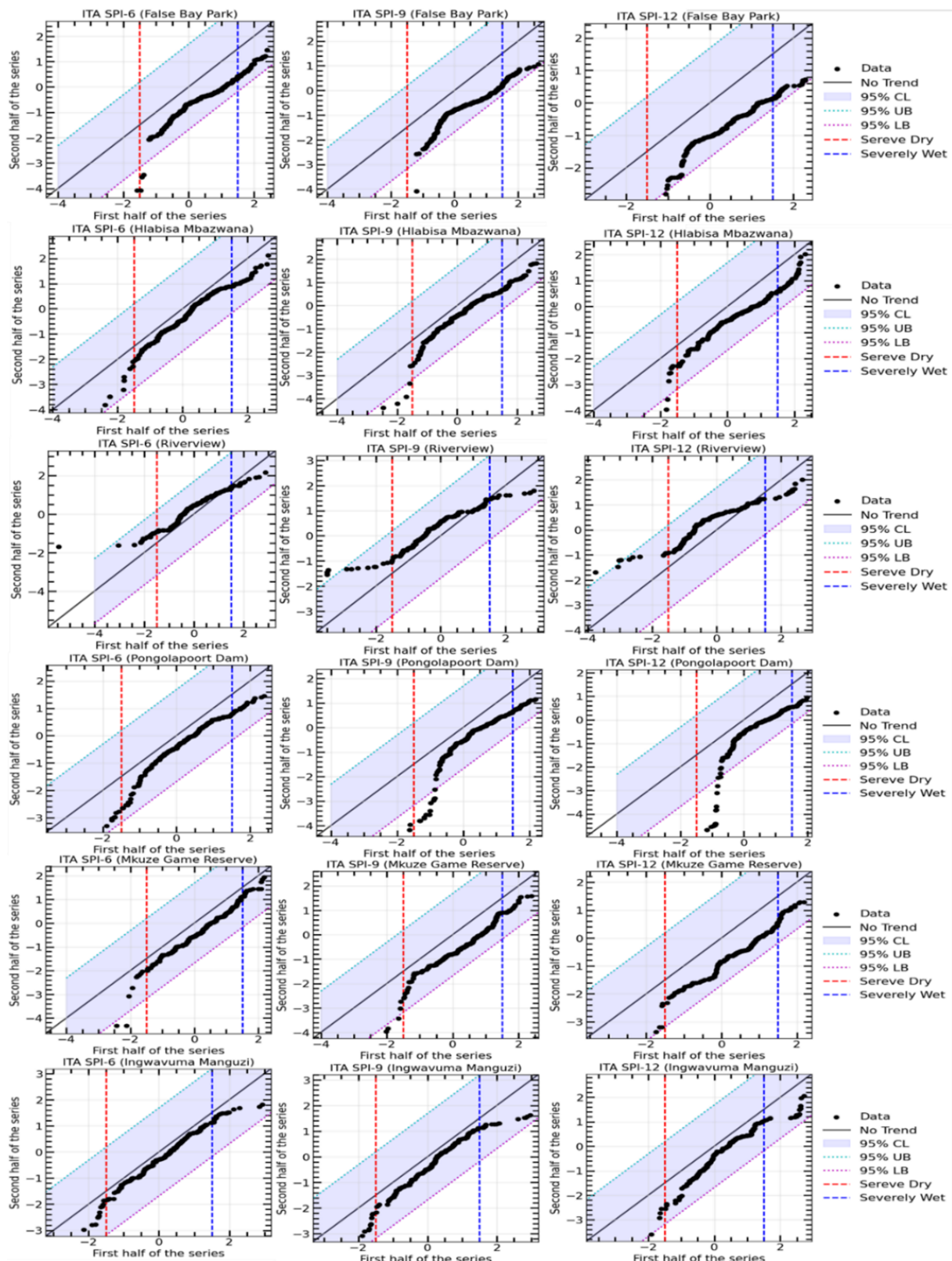
**Figure 8.** Results of Innovative trend analysis applied to different time scales values (SPI-6 (left), SPI-9 (middle), SPI-12 (right)). The blue shaded area represents the 95 % confidence level area. The red and blue vertical lines represent the severe drought and severely wet, respectively.

**Table 2.** Statistical summary of trend analysis for SPI-6, SPI-9, and SPI-12 using Mann-Kendall (MK) and Modified Mann-Kendall (MMK) tests.

| False Bay Park | | | |
| --- | --- | --- | --- |
| Variables | SPI-6 | SPI-9 | SPI-12 |
| $Z_{MK}$ | $-10.89$ | $-12.89$ | $-13.82$ |
| $p$-value$_{Mk}$ | $< 0.00$ | $< 0.00$ | $< 0.00$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| $Z_{MMK}$ | $-6.27$ | $-6.28$ | $-6.29$ |
| $p$-value$_{MMk}$ | $3.66 \times 10^{-10}$ | $3.35 \times 10^{-10}$ | $3.13 \times 10^{-10}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| **Hlabisa Mbazwana** | | | |
| $Z_{MK}$ | $-2.89$ | $-3.88$ | $-5.31$ |
| $p$-value$_{Mk}$ | $3.77 \times 10^{-3}$ | $3.05 \times 10^{-4}$ | $1.10 \times 10^{-7}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| $Z_{MMK}$ | $-2.26$ | $-2.12$ | $-2.20$ |
| $p$-value$_{MMk}$ | $2.39 \times 10^{-2}$ | $3.36 \times 10^{-2}$ | $2.78 \times 10^{-2}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| **Pongolapoort Dam** | | | |
| $Z_{MK}$ | $-7.19$ | $-8.74$ | $-9.83$ |
| $p$-value$_{Mk}$ | $6.12 \times 10^{-13}$ | $< 0.00$ | $< 0.00$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| $Z_{MMK}$ | $-8.22$ | $-5.44$ | $-6.51$ |
| $p$-value$_{MMk}$ | $2.22 \times 10^{-16}$ | $5.40 \times 10^{-8}$ | $7.41 \times 10^{-11}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| **Mkuze Game Reserve** | | | |
| $Z_{MK}$ | $-3.66$ | $-5.54$ | $-6.67$ |
| $p$-value$_{Mk}$ | $2.48 \times 10^{-4}$ | $2.99 \times 10^{-8}$ | $2.55 \times 10^{-11}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| $Z_{MMK}$ | $-2.44$ | $-2.79$ | $-2.22$ |
| $p$-value$_{MMk}$ | $1.46 \times 10^{-2}$ | $5.13 \times 10^{-3}$ | $2.64 \times 10^{-2}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| **Ingwavuma Manguzi** | | | |
| $Z_{MK}$ | $-2.38$ | $-3.72$ | $-4.92$ |
| $p$-value$_{Mk}$ | $1.72 \times 10^{-2}$ | $1.98 \times 10^{-4}$ | $8.72 \times 10^{-7}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| $Z_{MMK}$ | $-1.61$ | $-2.48$ | $-2.27$ |
| $p$-value$_{MMk}$ | $1.08 \times 10^{-1}$ | $1.31 \times 10^{-2}$ | $2.29 \times 10^{-2}$ |
| Decision (Trend$_{MK}$) | Decreasing | Decreasing | Decreasing |
| **Riverview** | | | |
| $Z_{MK}$ | $2.85$ | $3.84$ | $4.59$ |
| $p$-value$_{Mk}$ | $4.34 \times 10^{-3}$ | $1.25 \times 10^{-4}$ | $4.25 \times 10^{-6}$ |
| Decision (Trend$_{MK}$) | Increasing | Increasing | Increasing |
| $Z_{MMK}$ | $1.94$ | $2.16$ | $2.29$ |
| $p$-value$_{MMk}$ | $5.12 \times 10^{-2}$ | $3.07 \times 10^{-2}$ | $2.19 \times 10^{-2}$ |
| Decision (Trend$_{MK}$) | Increasing | Increasing | Increasing |

Step 1: Data Preprocessing Phase: To enhance the quality of the data and prepare it for decomposition, the original SPI time series undergo a data preprocessing phase:

– Savitzky–Golay Filter: This filter is applied to smooth the SPI data and preserves the essential shape and trends of the original time series while minimizing high-

**Table 3.** The results of the trend analysis for SPI-6, SPI-9, and SPI-12 obtained through a two-tailed test at a significance level of 5 % using ITA technique.

| False Bay Park | | | |
|---|---|---|---|
| Variables | SPI-6 | SPI-9 | SPI-12 |
| Slope | $-3.51 \times 10^{-3}$ | $-1.14 \times 10^{-3}$ | $-4.49 \times 10^{-3}$ |
| Indicator | $-20.08$ | $-20.12$ | $-20.07$ |
| $\pm$CI at 95 % | $\pm 9.24 \times 10^{-5}$ | $\pm 7.52 \times 10^{-5}$ | $\pm 6.82 \times 10^{-5}$ |
| Hlabisa Mbazwana | | | |
| Slope | $-1.68 \times 10^{-3}$ | $-2.31 \times 10^{-3}$ | $-1.86 \times 10^{-3}$ |
| Indicator | $-20.52$ | $-20.72$ | $-20.64$ |
| $\pm$CI at 95 % | $\pm 6.81 \times 10^{-5}$ | $\pm 9.35 \times 10^{-5}$ | $\pm 7.15 \times 10^{-5}$ |
| Pongolapoort Dam | | | |
| Slope | $2.26 \times 10^{-3}$ | $-2.88 \times 10^{-3}$ | $-3.34 \times 10^{-3}$ |
| Indicator | $-19.27$ | $-19.40$ | $-19.55$ |
| $\pm$CI at 95 % | $\pm 2.22 \times 10^{-5}$ | $\pm 3.62 \times 10^{-5}$ | $\pm 6.72 \times 10^{-5}$ |
| Mkuze Game Reserve | | | |
| Slope | $-2.00 \times 10^{-3}$ | $-3.04 \times 10^{-3}$ | $-3.80 \times 10^{-3}$ |
| Indicator | $-20.09$ | $-20.22$ | $-20.25$ |
| $\pm$CI at 95 % | $\pm 2.81 \times 10^{-3}$ | $\pm 4.67 \times 10^{-3}$ | $\pm 4.40 \times 10^{-3}$ |
| Ingwavuma Manguzi | | | |
| Slope | $-1.61 \times 10^{-3}$ | $-2.26 \times 10^{-3}$ | $-2.88 \times 10^{-3}$ |
| Indicator | $-21.96$ | $-21.05$ | $-20.77$ |
| $\pm$CI at 95 % | $\pm 6.81 \times 10^{-5}$ | $1.01 \pm \times 10^{-5}$ | $\pm 1.19 \times 10^{-5}$ |
| Riverview | | | |
| Slope | $1.69 \times 10^{-3}$ | $2.19 \times 10^{-3}$ | $2.37 \times 10^{-3}$ |
| Indicator | $22.54$ | $22.22$ | $21.86$ |
| $\pm$CI at 95 % | $\pm 1.54 \times 10^{-5}$ | $\pm 1.35 \times 10^{-5}$ | $\pm 1.56 \times 10^{-5}$ |

frequency noise. This step ensures that important signal patterns are retained during further processing. The smoothed signal becomes the input signal for decomposition technique.

– CEEMDAN Parameter Settings: CEEMDAN is used to break the smoothed signal into several IMFs and a residual component. Before decomposition, the necessary parameters for CEEMDAN are configured. These parameters control the number of realizations, noise amplitude, and stopping criteria for decomposition.

Step 2: Model Development Phase: Each IMF, including the residual, is independently modelled using a hybrid ARIMA–LSTM approach. This process involves several steps:

a. Data Partitioning

– The data for each IMF is split into: Training set (80 %) and Testing set (20 %). This split ensures

that model learning and evaluation are based on separate subsets to avoid overfitting.

b. Normalization

– Prior to model training, the data is normalized using Min-Max normalization to ensure that input features fall within a similar scale, which improves training stability and convergence speed.

c. Modelling Each IMF with ARIMA–LSTM

– The two models are integrated so that both linear (ARIMA) and nonlinear (LSTM) dependencies within each IMF are effectively captured. The modelling process follows the algorithm shown in Fig. 4.

d. Feature Selection and Hyperparameter Tuning

– The performance of ARIMA and LSTM models heavily depends on the feature selection and
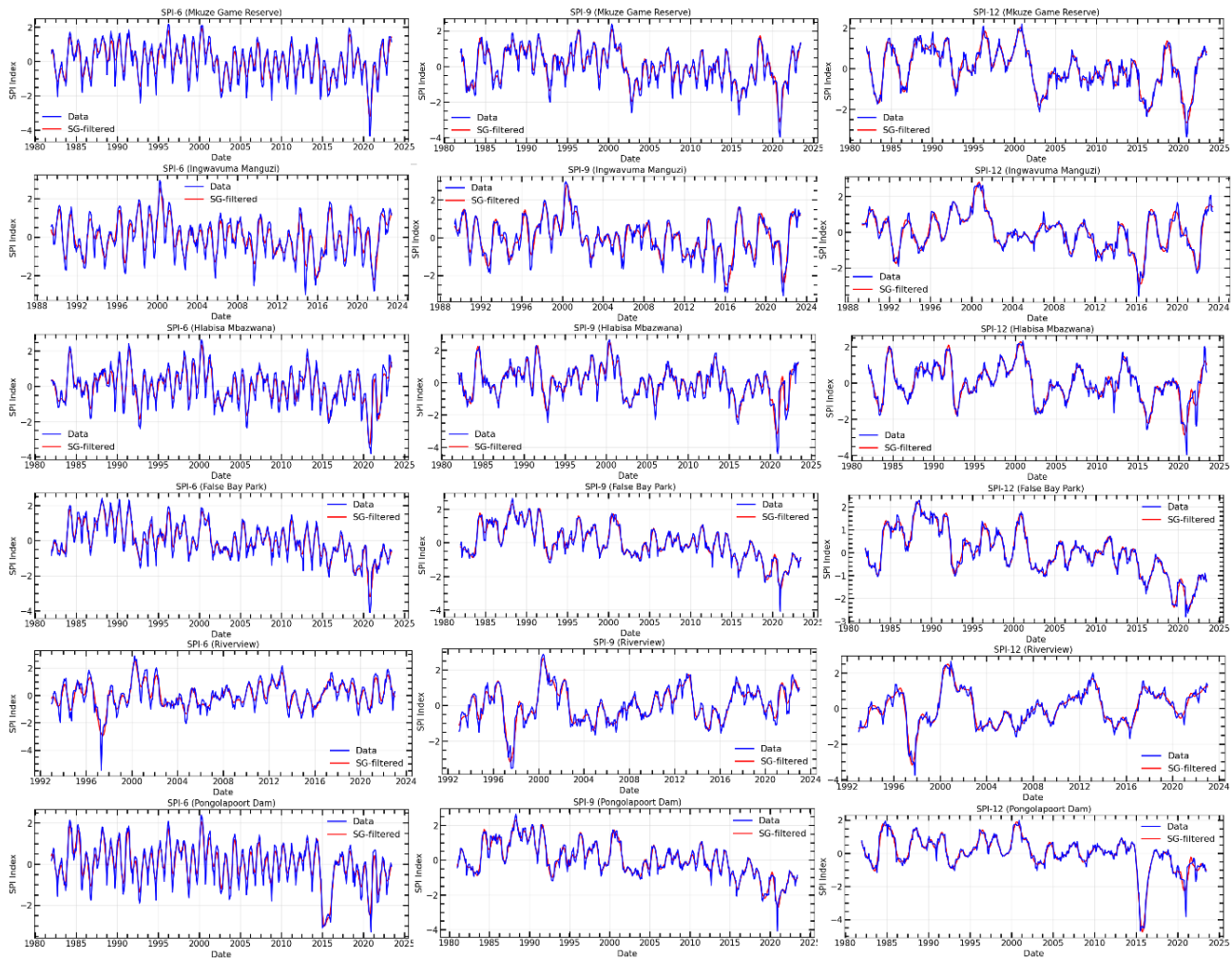
**Figure 9.** SPI signals smoothed by Savitzky-Golay (SG).

hyperparameters. The auto_arima() function and Bayesian Optimization were used to automate and optimize the search for best-performing hyperparameters for the ARIMA-LSTM model by evaluating model performance over a probabilistic space.

e. Model Training

– Each IMF is trained individually using the selected features and optimized hyperparameters, resulting in a trained model for each component.

Step 3: Forecast Reconstruction Phase

– After training, each IMF is forecasted individually. The final forecasted SPI value is obtained by summing the predictions of all individual IMFs, including the residual component:

$$\hat{\text{SPI}}(t) = \sum_{i=1}^{n} I\hat{M}F_i(t) + \text{Res}_t$$

This additive reconstruction ensures that the original structure and dynamics of the SPI series are preserved in the forecast, improving overall accuracy.

Step 4: Model Evaluation Phase

The reconstructed SPI prediction is then evaluated using multiple performance metrics: RMSE, DS, and coefficient of determination. The Taylor diagram is also utilised to evaluate the model performance. These metrics help quantify the predictive accuracy and reliability of the hybrid framework.

## 2.11   Performance Evaluation

To establish the predictive superiority of the SG-CEEMDAN-ARIMA-LSTM model, a comparison was conducted against other models, including ARIMA, LSTM, ARIMA-LSTM, SG-ARIMA-LSTM, and CEEMDAN-ARIMA-LSTM models. The performance of the proposed hybrid-based model is evaluated using three indicators namely, root mean square error (RMSE), coefficient of
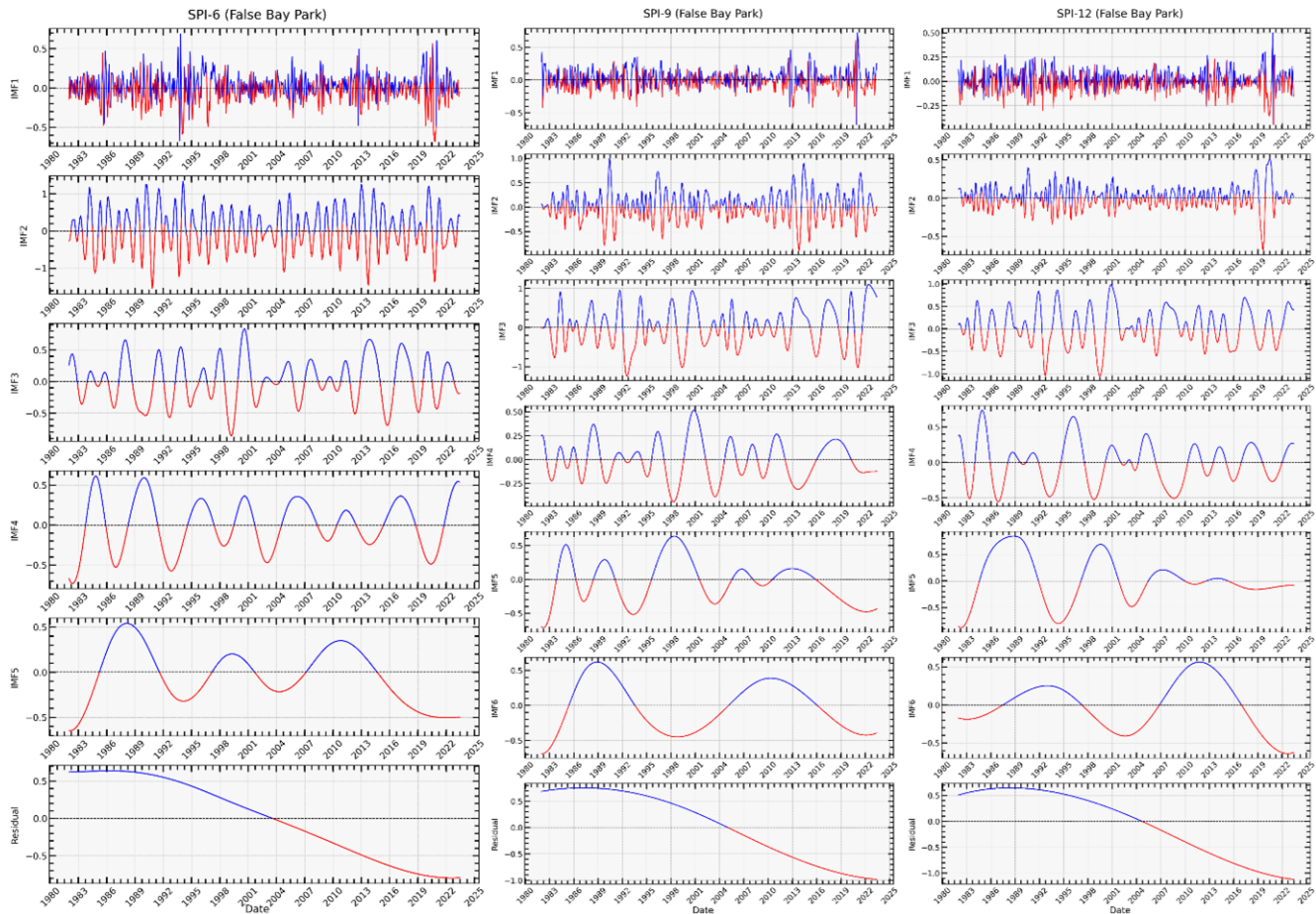
**Figure 10.** Decomposition of Smoothed SPI-6, SPI-9 and SPI-12 Index Using CEEMDAN: Each IMF represents different frequency components of the SPI index, from high-frequency oscillations (IMF1) to low-frequency trends (IMF5), showing the variability in precipitation patterns over the years from 1980 to 2023.

determination $(R^2)$ and directional symmetry (DS). The high value of $R^2$ and DS reflects the better performance of the forecasting model while the lower the value of RMSE illustrates better forecasting performance.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left( y_i - \hat{y}_{\text{avg}} \right)^2} \tag{43}$$

$$R^2 = \frac{\left[ \sum_{i=1}^{n} \left( y_i - y_{\text{avg}} \right) \left( \hat{y}_i - \hat{y}_{\text{avg}} \right) \right]^2}{\sum_{i=1}^{n} \left( y_i - y_{\text{avg}} \right)^2 \sum_{i=1}^{n} \left( \hat{y}_i - \hat{y}_{\text{avg}} \right)^2} \tag{44}$$

$$\text{DS} = \frac{100}{n-1} \sum_{i=2}^{n} d_i \tag{45}$$

where

$$d_i = \begin{cases} 1, & \left( y_i - y_{i-1} \right) \left( \hat{y}_i - \hat{y}_{i-1} \right) > 0 \\ 0, & \text{otherwise} \end{cases} \tag{46}$$

$n$ is number of data points, $y_i$ and $\hat{y}_i$ observed and forecasted, respectively. $y_{\text{avg}}$ and $\hat{y}_{\text{avg}}$ an average of the actual and forecasted values, respectively. Furthermore, this study conducts a qualitative evaluation of the prediction model's performance using a Taylor diagram (Taylor, 2001). The Taylor diagram offers a statistical evaluation of the degree of agreement between the models in terms of their SD, RMSE, and $R^2$, while providing a concise summary of the correspondence between predicted and observed values. The differences in DS, RMSE, and $R^2$ values among the prediction models are depicted as individual points on a two-dimensional plot within the Taylor diagram. This diagram, though it follows a common structure, proves especially valuable when evaluating intricate models.

**Table 4.** ADF Test Results for SPI Values (SPI-6, SPI-9, SPI-12) at Different Stations.

| Station Name | SPI | ADF Statistic | *p*-value | Critical Value (5 %) |
|---|---|---|---|---|
| False Bay Park | SPI-6 | −2.1926 | 0.2089 | −2.8925 |
| | SPI-9 | −3.2142 | 0.0192 | −2.8915 |
| | SPI-12 | −1.4829 | 0.5419 | −2.8949 |
| Hlabisa Mbazwana | SPI-6 | −1.9314 | 0.3175 | −2.8925 |
| | SPI-9 | −1.5629 | 0.5022 | −2.8939 |
| | SPI-12 | −1.1867 | 0.6793 | −2.8946 |
| Pongolapoort Dam | SPI-6 | −2.8759 | 0.0482 | −2.8925 |
| | SPI-9 | −2.7909 | 0.0596 | −2.8909 |
| | SPI-12 | −2.1864 | 0.2112 | −2.8909 |
| Mkuze Game Reserve | SPI-6 | −3.1136 | 0.0256 | −2.8949 |
| | SPI-9 | −1.6134 | 0.4762 | −2.8939 |
| | SPI-12 | −2.5689 | 0.0996 | −2.8949 |
| Ingwavuma Manguzi | SPI-6 | −2.1418 | 0.2281 | −2.8994 |
| | SPI-9 | −3.6158 | 0.0055 | −2.9026 |
| | SPI-12 | −1.9049 | 0.3298 | −2.9026 |
| Riverview | SPI-6 | −1.7509 | 0.4051 | −2.9051 |
| | SPI-9 | −1.1840 | 0.6804 | −2.9079 |
| | SPI-12 | −2.0298 | 0.2737 | −2.9015 |

## 3 Results and Discussion

### 3.1 Rainfall Data Series

Figure 6 illustrates the daily and monthly cumulative precipitation data recorded at the uMkhanyakude district meteorological stations in KwaZulu-Natal province, South Africa, from the early 1980s to 2023. The data comprising 20 % was employed for prediction, whereas the data representing 80 % was applied for training. The SPI was computed utilising rainfall data from meteorological stations in the uMkhanyakude district, which provide sufficiently extensive records and a consistent structure (Hırca et al., 2022).

### 3.2 SPI Time Series and Trend Analysis

This study SPI values for the 6-, 9-, and 12-month intervals were computed using the monthly mean time series shown in Fig. 6. Figure 7 illustrates the time series of the SPI calculated for the 6-month (SPI-6), 9-month (SPI-9), and 12-month (SPI-12) intervals. All SPIs (SPI-6, SPI-9, and SPI-12) demonstrate numerous occurrences of moderate to severe droughts in the studied area. A significant drought episode was reported from late 2004 to 2009. Moreover, SPI-12 exhibits a persistent drought spell that commenced between 2014 and 2016, resulting in a decline in water supply conditions in the region (Bukhosini and Moyo, 2023). The statistics across all timelines indicate a troubling trend of extended and intense drought conditions in recent years. This under-

scores the pressing necessity for efficient water management and drought readiness in the area. Initially, we assess the trend throughout the research area employing nonparametric techniques. The ensuing conclusions will be obtained via advanced trend analysis methods employed to investigate SPI trends.

Figure 8 illustrates the regional outcomes of the ITA methodology used on the 6-, 9-, and 12-month SPI series to ascertain the potential meteorological drought trend in the uMkhanyakude district. Figure 8 includes two vertical bands to elucidate the potential trends of arid and humid conditions: a red band indicating the drought threshold (SPI = −1.5) and a blue band denoting the wet threshold (SPI = 1.5). The zone between the two bands signifies normal conditions, hence facilitating the depiction of both low and high SPI trends using the ITA methodology. Each plot compares the first and second halves of the data series to identify trends.

In general, both Fig. 8 and Table 3 show that all stations, except Riverview, indicate a downward trend for all time scales, in terms of the ITA. For example, the ITA results obtained using 6-month SPI values exhibit a slightly decreasing trend in precipitation, moving toward the upper right quadrant, indicating the detection of drier conditions over the 6-month timescale. Some points approach the severely wet threshold but do not cross it, indicating that there were no extreme wet periods, though some drier periods are evident near the severe dry line. The ITA results obtained using 9-month SPI values show a more pronounced decreasing trend, indicating a relatively weaker increase in wet conditions over a 9-month timescale. Several points approach the severe dry threshold, but the data remains mostly within the 95 % confidence bounds, indicating moderate variability in precipitation trends. On the other hand, the SPI-12 plot demonstrates a noticeable decreasing trend toward dryness, as many points fall below the no-trend line and approach the severe dry region. Riverview indicates the increasing trend across all time scales. The increasing distance between the black dots and the no-trend line highlights a shift toward drier conditions in the second half of the series. In general, the analysis suggests a gradual increase in precipitation for shorter periods (SPI-6), moderate upward trends for medium-term periods (SPI-9), and a more substantial shift toward dry conditions over longer periods (SPI-12) for Riverview. The variability is evident, but a clear progression toward drier conditions is evident, particularly in the SPI-12 plot. This observation could be indicative of changing precipitation patterns, which is crucial for understanding drought risk and informing water resource management strategies.

Table 2 presents the results of the Mann-Kendall (MK) and Modified Mann-Kendall (MMK) trend tests for the Standardized Precipitation Index (SPI) over 6-month (SPI-6), 9-month (SPI-9), and 12-month (SPI-12) periods. The results indicate that across five stations all time scales both MK and MMK methods showed significant decreasing trend with negative Z-score values. For example, False Bay Park,
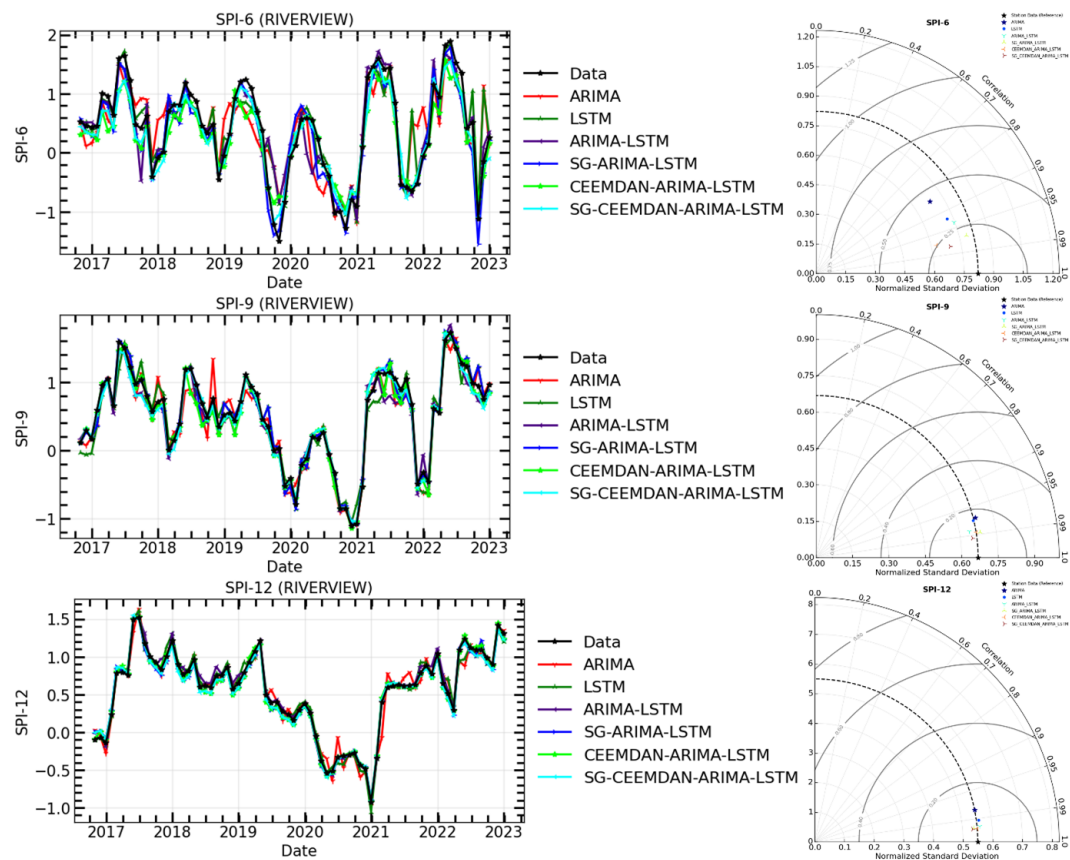
**Figure 11.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of Riverview meteorological station.

Z_MK are $Z_{\text{SPI-6}} = -10.89$, $Z_{\text{SPI-9}} = -12.89$, $Z_{\text{SPI-12}} = -13.82$ and Z_MMK are $Z_{\text{SPI-6}} = -6.27$, $Z_{\text{SPI-9}} = -6.28$, $Z_{\text{SPI-12}} = -6.29$. The $p$-values of MK and MMK show the significance of the trends, with values way below 0.05 confirming statistically significant trends. In all cases except Riverview, the $p$-values are extremely low ($\ll 0.05$), indicating strong evidence of significant decreasing trends in precipitation for all SPI periods. Both the MK and MMK tests confirm decreasing trends across all time scales, with the Z_MK and Z_MMK values becoming more negative as the SPI period increases, reflecting an intensifying downward trend over longer periods (from SPI-6 to SPI-12). For Riverview station, the results indicate an increasing trend with positive Z-score values, i.e. Z_MK are $Z_{\text{SPI-6}} = 2.85$, $Z_{\text{SPI-9}} = 3.84$, $Z_{\text{SPI-12}} = 4.59$ and Z_MMK are $Z_{\text{SPI-6}} = 1.19$, $Z_{\text{SPI-9}} = 2.16$, $Z_{\text{SPI-12}} = 2.29$. In general, all these results are consistent with those shown using the ITA (see Table 3). The Riverview station experience increasing trend because it is located closer to the coast, hence it is influenced by a combination of geographic, oceanic and climatic factors. For an example, this station could be influenced by the Agulhas Current, which flows southwards along the east coast of South Africa, bringing warm, moist air from the Indian Ocean, and

thus enhancing evaporation that brings constant availability of moisture in the atmosphere.

## 3.3 SPI Time Series Forecasting Results

The study proposes a hybrid model that applies the Savitzky-Golay (SG) filter to process raw SPI data, thereby reducing noise and enhancing forecasting analysis. To demonstrate the effectiveness of the SG filter, appropriate parameters such as window size and polynomial order were selected through trial and error using data from the study sites (Sibiya et al., 2024). A window size of 21 and a polynomial order of 5 were chosen for smoothing. Figure 9 shows how the SG filter effectively tracks the general trend while preserving the shape of peaks and minimizing noise. This filter was applied to different time scales of the SPI time series. It autonomously calibrates according to peak distribution, exhibiting optimal performance, particularly with asymmetric peaks, while preserving peak height integrity. The application of the SG filter effectively mitigates short-term fluctuations and eliminates noise from the time series, resulting in cleaner data, thereby enhancing the reliability of the subsequent decomposition process. By reducing noise, decomposition techniques can
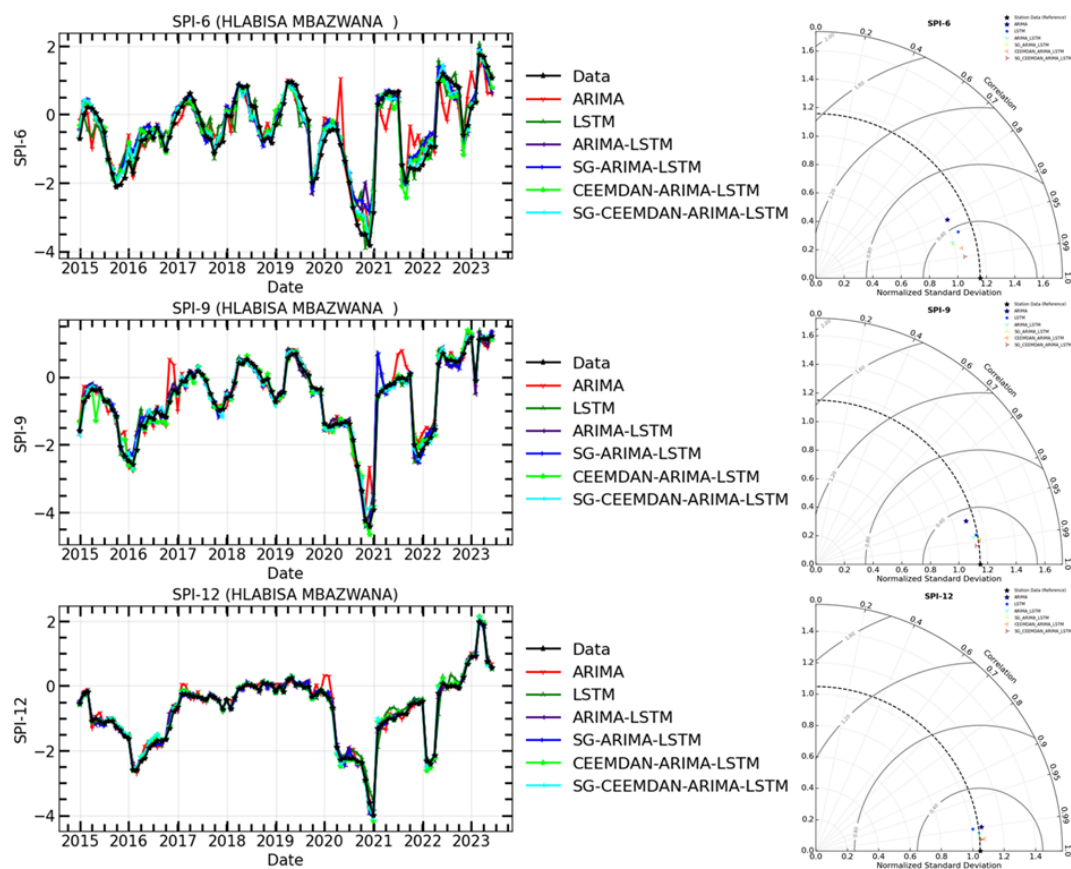
**Figure 12.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of Hlabisa Mbazwana meteorological station.

more accurately capture the authentic underlying patterns and components within the data.

After applying a Savitzky-Golay filter to the series, the CEEMDAN algorithm decomposes the filtered SPI series into six subseries with different amplitudes and frequencies. The results from the False Bay Park station are utilized here as an illustration to prevent repetition. In these results, the decomposed set of time series consists of five IMF components and a residual component, as shown in Fig. 10 (for all time scales). During the decomposition process, white Gaussian noise is added to create noisy signals. The original sequence exhibits high nonlinearity and nonstationarity, with the frequency of the IMF components gradually decreasing. Figure 10 depicts this gradual decrease in frequency as the order of the IMF components increases. As each IMF is further decomposed, it becomes less volatile and cyclical, which aligns with the characteristics of the decomposed IMF. Therefore, by predicting each IMF and the residual, the forecast precision can be enhanced. A forecasting model is then constructed for each component, and the prediction results are obtained by summing up the outputs of all predicted components.

In predictive modeling, this study employed Bayesian optimization for hyperparameter tuning because of its effectiveness in improving model performance for complex, black-box, and non-differentiable functions. The hyperparameter configuration space comprises an $n$-dimensional functional space that encompasses all possible combinations of hyperparameters for the specified model. The benchmark analysis began with the ARIMA model, using the Box–Jenkins methodology. This process started with an assessment of stationarity through the augmented Dickey–Fuller (ADF) test. The series showed $p$-values exceeding the 5 % significance threshold, indicating non-stationarity (see Table 4). As a result, differencing was applied to achieve stationarity. This study employed a stepwise approach using the auto_arima( ) function within the ARIMA framework to identify the optimal parameters (see Table 5). Table 6 delineates the hyperparameter search space employed for tuning the LSTM model utilizing a Bayesian optimization approach. Each hyperparameter is presented alongside its respective range or selected value, which delineates the parameters within which the Bayesian search investigated optimal configurations.

The models in Table 7 were compared for their prediction ability before and after time series decomposition in this
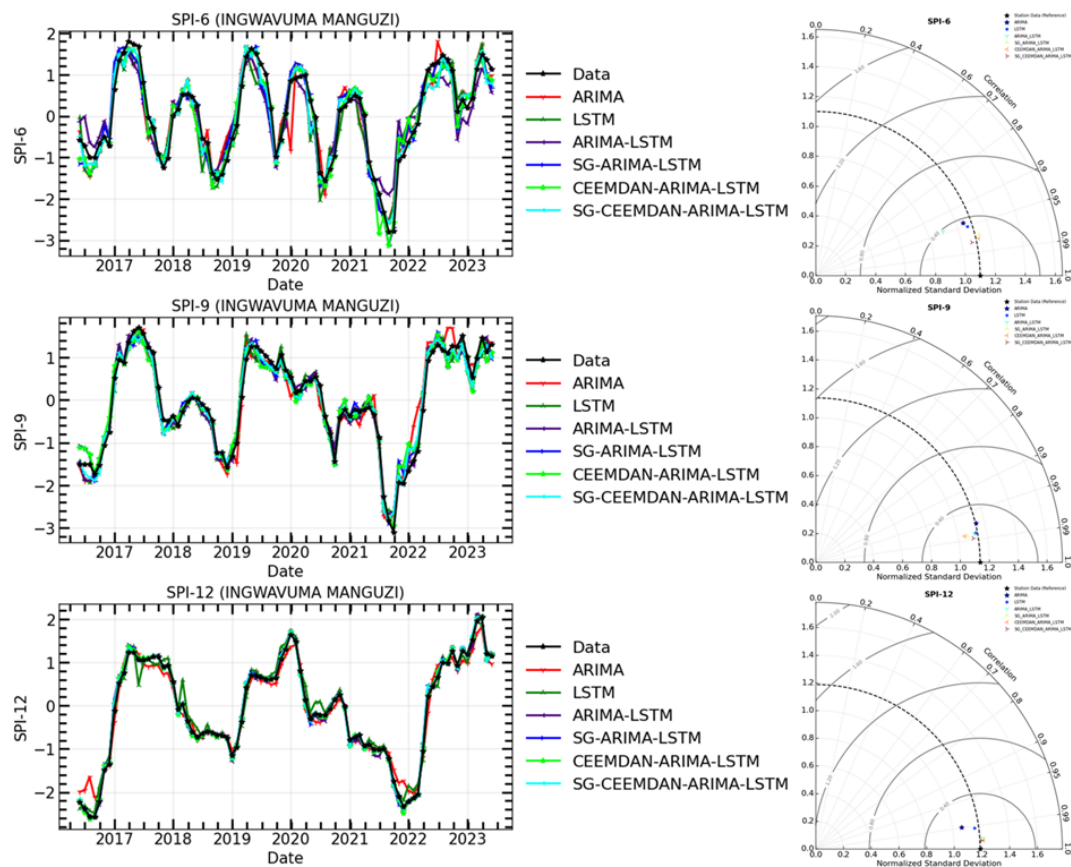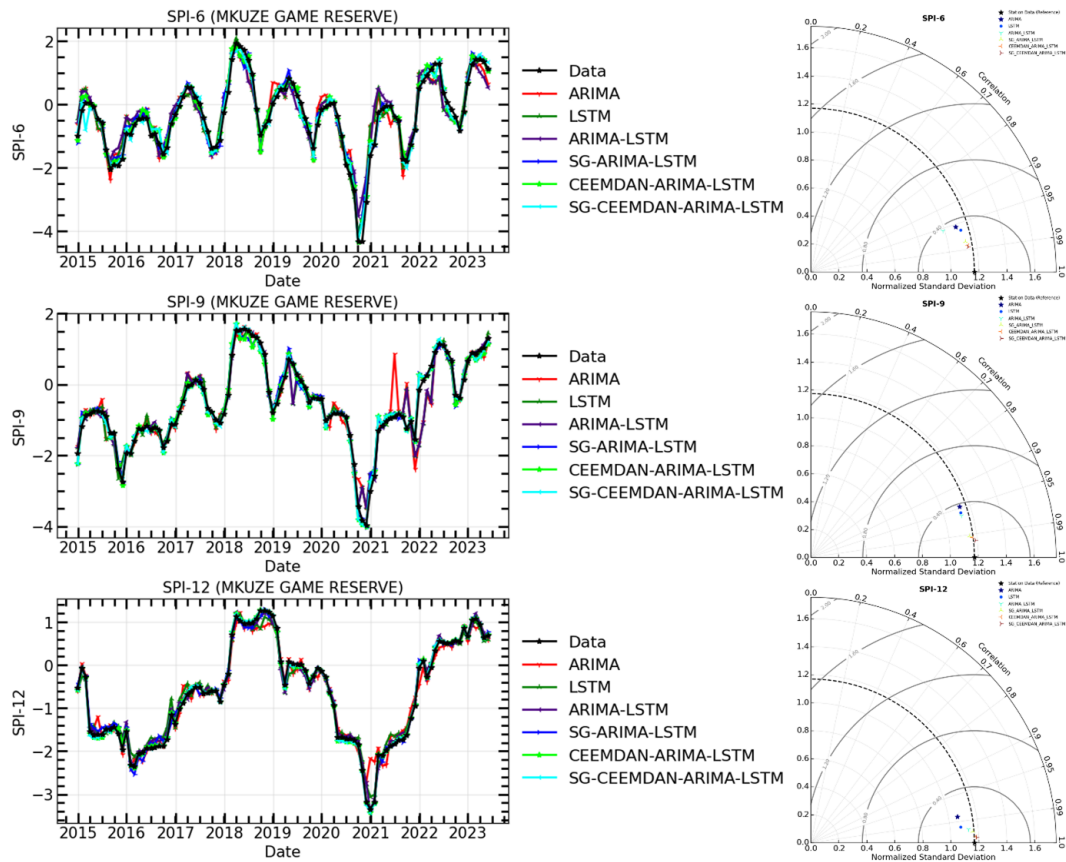
**Figure 13.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of Ingwavuma Manguzi meteorological station.

**Table 5.** Accuracy criteria for different model parameters of the ARIMA model applied in SPI-6, SPI-9 and SPI-12 at different meteorological stations of uMkhanyakude district.

| Station Name | SPI-6 | | SPI-9 | | SPI-12 | |
|---|---|---|---|---|---|---|
| | Model | AIC | Model | AIC | Model | AIC |
| False Bay Park | ARIMA(5,0,3) | 517.757 | ARIMA(3,1,1) | 333.328 | ARIMA(1,1,0) | 183.988 |
| Hlabisa Mbazwana | ARIMA(5,1,5) | 322.514 | ARIMA(3,0,5) | 248.815 | ARIMA(2,1,2) | 152.295 |
| Pongolapoort Dam | ARIMA(4,1,3) | 438.230 | ARIMA(3,1,2) | 350.618 | ARIMA(1,1,0) | 254.076 |
| Mkuze Game Reserve | ARIMA(4,1,2) | 432.320 | ARIMA(3,0,3) | 330.540 | ARIMA(0,1,1) | 164.170 |
| Ingwavuma Manguzi | ARIMA(4,0,5) | 417.071 | ARIMA(3,1,1) | 350.196 | ARIMA(0,1,1) | 153.087 |
| Riverview | ARIMA(4,1,5) | 435.687 | ARIMA(3,1,0) | 365.509 | ARIMA(2,1,1) | 168.812 |

research. The objective was to determine if smoothing and decomposing time series improve the model's prediction performance. Figures 11–16 show a comparison of the various models' prediction outcomes using the Taylor diagram. In general, all the models accurately replicate the original SPI time series at all timescales (refer to Figs. 11–16) in terms of the time series plot. However, the SG-CEEMDAN-ARIMA-LSTM model (shown in red) appears to have the closest fit to the data, displaying superior accuracy across different phases, particularly in extreme values. Nonetheless, the

hybrid models (SG-ARIMA-LSTM, CEEMDAN-ARIMA-LSTM, and SG-CEEMDAN-ARIMA-LSTM) show better precision in capturing peaks, rapid transitions and troughs compared to the standalone LSTM or ARIMA models. Table 7 displays an assessment of the predictive performance metrics of several models utilising RMSE, $R^2$, and DS. As the period extends, the RMSE values decrease; however, the DS and cap $R$-squared values typically enhance (see Table 7). This indicates that the models' predictive accuracy progressively enhances with an extended duration, reach-

**Table 6.** Hyperparameter ranges in LSTM–Bayesian search Method.

| Hyperparameters | Values | Hyperparameters | Values |
|---|---|---|---|
| Number of LSTM units | (32, 256) | Activation function | (ReLu, Sigmoid, Tanh,) |
| Number of LSTM hidden size | (32, 256) | Optimizer | Adam |
| Batch size | (16,128) | Loss function | Mean Square error |
| Epoch | (50,300) | Dropout | (0.05, 0.1) |
| LSTM learning rate | (0.0001, 0.001) | Regularization | Early stopping |



**Figure 14.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of Mkuze Game Reserve meteorological station.

ing its highest point at the 12-month interval. In terms of RMSE, the SG-CEEMDAN-ARIMA-LSTM model outperforms the others, exhibiting the lowest error values across all indices. For example, Riverview station, 0.2165 for SPI-6, 0.0921 for SPI-9, and 0.0566 for SPI-12. This indicates that this model has the smallest prediction error, making it the most accurate in terms of error reduction. Concerning $R^2$, which measures how well the model explains the variance in the data, SG-CEEMDAN-ARIMA-LSTM again leads with the highest values: 0.9602 for SPI-6, 0.9846 for SPI-9, and 0.9939 for SPI-12. This shows that the model provides the best fit to the data. The CEEMDAN-ARIMA-LSTM model is the second-best performer, also exhibiting impressive re-

sults, particularly in $R^2$, where it achieves higher values of 0.9483 for SPI-6, 0.9751 for SPI-9, and 0.9933 for SPI-12. The SG-ARIMA-LSTM model is the third-best hybrid performer, with RMSE values of 0.2262 for SPI-6, 0.1051 for SPI-9, and 0.05639 for SPI-12. The SG-ARIMA-LSTM model is the third-best performer, also exhibiting impressive results, particularly in $R^2$, where it achieves higher values of 0.9392 for SPI-6, 0.9763 for SPI-9, and 0.9904 for SPI-12. The SG-ARIMA-LSTM model is the third-best hybrid performer, with RMSE values of 0.2597 for SPI-6, 0.1157 for SPI-9, and 0.0567 for SPI-12. In general, these results highlight the efficacy of hybrid models, particularly those incorporating SG and CEEMDAN processes, in improving
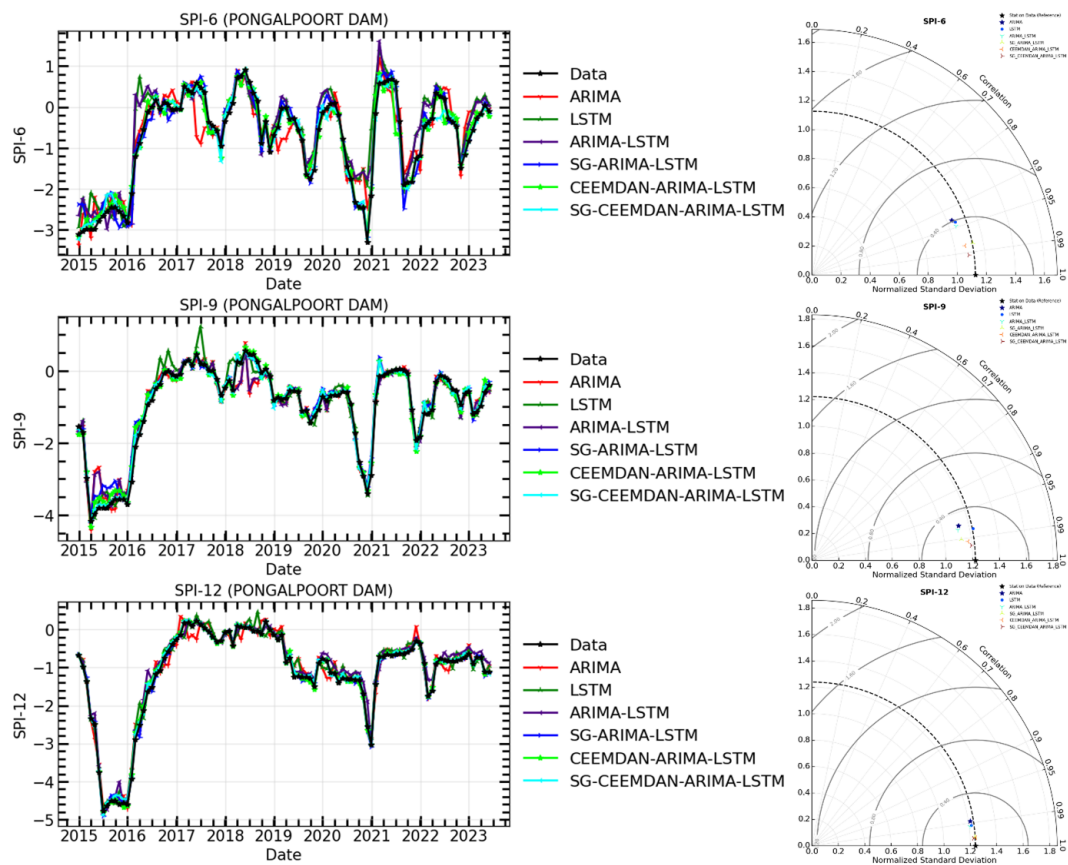
**Figure 15.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of Pongolapoort dam meteorological station.

predictive accuracy across multiple timescales of SPI, particularly for the SG-CEEMDAN-ARIMA-LSTM model. These results are consistent with the Taylor diagram (see Figs. 11–16), which indicates a significant improvement in prediction accuracy after incorporating the SG and CEEMDAN signal decomposition technique as the hybrid model exhibits superior performance in terms of prediction accuracy across all timescales, surpassing other models. This suggests that the inclusion of these techniques enhances the models' ability to capture both short-term and long-term dependencies, thus making them more robust for drought prediction purposes. Therefore, this hybrid model appears to be the most effective for drought prediction in this analysis. These findings highlight the superiority of the proposed hybrid model in enhancing drought prediction accuracy compared to standalone approaches.

## 4 Discussion

In this study, we utilized the Mann-Kendall and Modified Mann-Kendall tests to determine the drought trend index in meteorological variables within the basin. The MK and MMK trend methods showed a significant decrease in all SPI time scales based on rainfall data from five stations; however, the district, except for the Riverview station, showed an increasing trend in the uMkhanyakude district. The study's findings align with prior research by Kganvago et al. (2021) and Ngwenya et al. (2024). Ngwenya et al. (2024) conducted a study using the Mann-Kendall test to assess the SPI values at a 5 % significance level, revealing sustained drought conditions in the Western Cape region. Kganvago et al. (2021) indicated a notable decline in drought conditions in the Western Cape area of South Africa. We have also employed the ITA, which enhances the MK and MMK tests in identifying trends, and the results underscore the importance of comprehending drought conditions. The findings of our analysis validate previous research by Naik and Abiodun (2020), highlighting the need to conduct trend studies on drought indicators to investigate the impacts of climate change. The study highlights the crucial role of SPI as a primary variable in monitoring and forecasting droughts in the region, and its potential to mitigate the adverse impacts of droughts and water scarcity in the uMkhanyakude district in the future. The objective was to determine if the model's predictive performance is enhanced by smoothing and deconstructing time series data.

**Table 7.** Performance measures for the comparison of observed and forecasted data of the models for SPI-6, SPI-9 and SPI-12 across various lead times using statistical criteria.

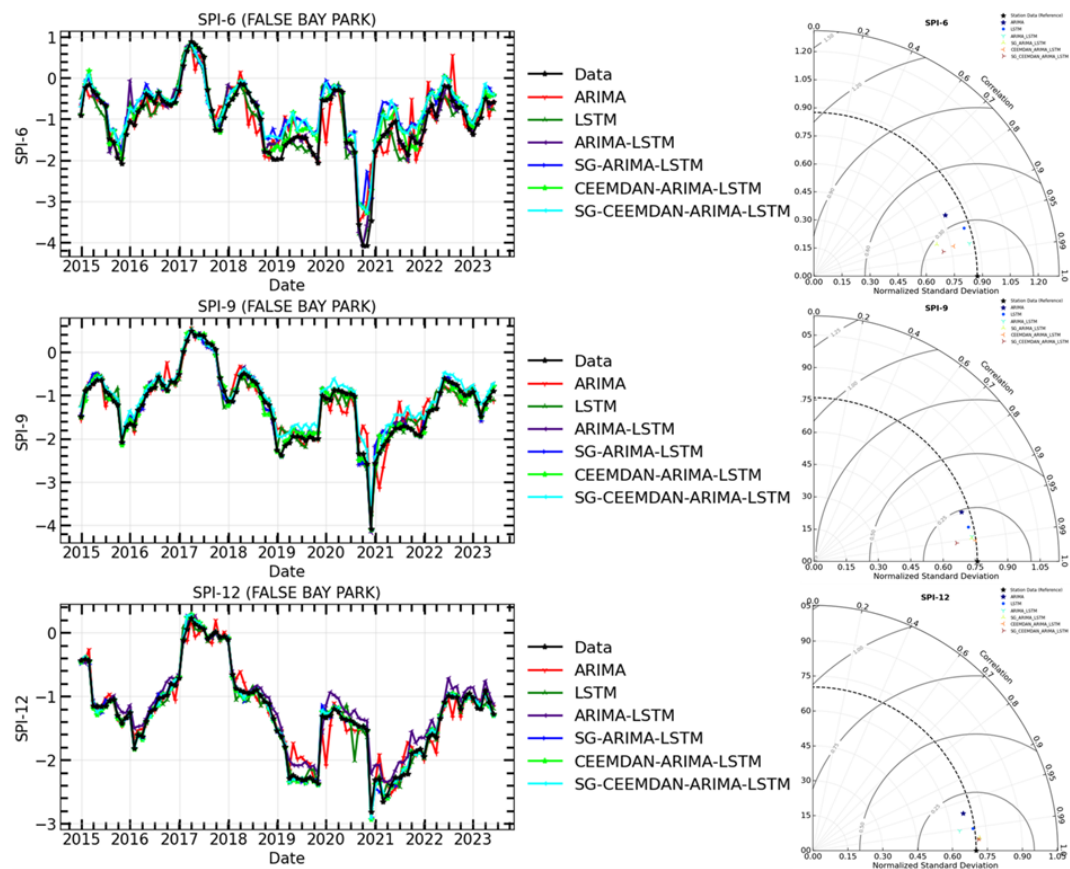| False Bay Park | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | SPI-6 | | | SPI-9 | | | SPI-12 | | |
| | RMSE | $R^2$ | DS | RMSE | $R^2$ | DS | RMSE | $R^2$ | DS |
| ARIMA | 0.3504 | 0.8435 | 0.8426 | 0.2431 | 0.8976 | 0.8525 | 0.1689 | 0.9421 | 0.8426 |
| LSTM | 0.3128 | 0.9111 | 0.8327 | 0.2416 | 0.9521 | 0.8723 | 0.1626 | 0.9821 | 0.8519 |
| ARIMA-LSTM | 0.2476 | 0.9194 | 0.8327 | 0.1650 | 0.9531 | 0.8723 | 0.0507 | 0.9952 | 0.9009 |
| SG-ARIMA-LSTM | 0.2056 | 0.9458 | 0.8030 | 0.1348 | 0.9687 | 0.8218 | 0.0571 | 0.9940 | 0.9009 |
| C-A-L | 0.2182 | 0.9375 | 0.8713 | 0.0978 | 0.9834 | 0.8218 | 0.0496 | 0.9953 | 0.8911 |
| SG-C-A-L | 0.1835 | 0.9650 | 0.8416 | 0.1631 | 0.9836 | 0.8317 | 0.0349 | 0.9957 | 0.8941 |
| **Mkuze Game Reserve** | | | | | | | | | |
| ARIMA | 0.3752 | 0.8642 | 0.8419 | 0.3475 | 0.8957 | 0.8792 | 0.2202 | 0.9697 | 0.8730 |
| LSTM | 0.3474 | 0.9121 | 0.8822 | 0.3354 | 0.9178 | 0.8030 | 0.1523 | 0.9890 | 0.8733 |
| ARIMA-LSTM | 0.3160 | 0.9273 | 0.8416 | 0.1561 | 0.9823 | 0.8218 | 0.1079 | 0.9926 | 0.8730 |
| SG-ARIMA-LSTM | 0.2307 | 0.9624 | 0.8515 | 0.1548 | 0.9825 | 0.8317 | 0.08252 | 0.9951 | 0.8019 |
| C-A-L | 0.1969 | 0.9726 | 0.8317 | 0.1430 | 0.9850 | 0.8515 | 0.04497 | 0.9986 | 0.9208 |
| SG-C-A-L | 0.1818 | 0.9742 | 0.8515 | 0.1232 | 0.9892 | 0.8617 | 0.04217 | 0.9990 | 0.9208 |
| **Pongolapoort Dam** | | | | | | | | | |
| ARIMA | 0.4470 | 0.8797 | 0.8624 | 0.2993 | 0.9668 | 0.8119 | 0.1918 | 0.9763 | 0.8733 |
| LSTM | 0.4470 | 0.8962 | 0.8732 | 0.2873 | 0.9467 | 0.8238 | 0.1824 | 0.9851 | 0.8829 |
| ARIMA-LSTM | 0.4121 | 0.8969 | 0.8822 | 0.2599 | 0.9588 | 0.8921 | 0.1638 | 0.9862 | 0.8432 |
| SG-ARIMA-LSTM | 0.2224 | 0.9617 | 0.8019 | 0.2064 | 0.9803 | 0.8515 | 0.0686 | 0.9969 | 0.8119 |
| C-A-L | 0.2132 | 0.9649 | 0.8822 | 0.1572 | 0.9850 | 0.8218 | 0.0639 | 0.9975 | 0.8019 |
| SG-C-A-L | 0.1453 | 0.9839 | 0.8824 | 0.1429 | 0.9858 | 0.8911 | 0.0635 | 0.9978 | 0.8921 |
| **Hlabisa Mbazwana** | | | | | | | | | |
| ARIMA | 0.4704 | 0.8347 | 0.8624 | 0.4234 | 0.8698 | 0.8921 | 0.2321 | 0.9556 | 0.8142 |
| LSTM | 0.3617 | 0.9041 | 0.8327 | 0.2163 | 0.9672 | 0.8119 | 0.1566 | 0.9806 | 0.8317 |
| ARIMA-LSTM | 0.3269 | 0.9369 | 0.8515 | 0.2139 | 0.9677 | 0.8218 | 0.1457 | 0.9813 | 0.8426 |
| SG-ARIMA-LSTM | 0.3011 | 0.9355 | 0.8416 | 0.1829 | 0.9747 | 0.8317 | 0.08540 | 0.9935 | 0.8218 |
| C-A-L | 0.2497 | 0.9592 | 0.8218 | 0.1662 | 0.9792 | 0.8218 | 0.0825 | 0.9949 | 0.9009 |
| SG-C-A-L | 0.1921 | 0.9795 | 0.8614 | 0.1332 | 0.9866 | 0.8218 | 0.07416 | 0.9952 | 0.9029 |
| **Ingwavuma Manguzi** | | | | | | | | | |
| ARIMA | 0.4123 | 0.8716 | 0.8571 | 0.2706 | 0.9442 | 0.8750 | 0.2052 | 0.9784 | 0.8619 |
| LSTM | 0.3843 | 0.8931 | 0.8738 | 0.2524 | 0.2524 | 0.8691 | 0.1614 | 0.9828 | 0.8095 |
| ARIMA-LSTM | 0.3458 | 0.9044 | 0.8095 | 0.2428 | 0.9695 | 0.8541 | 0.8541 | 0.9847 | 0.8215 |
| SG-ARIMA-LSTM | 0.2767 | 0.9397 | 0.8076 | 0.2001 | 0.9724 | 0.8809 | 0.0815 | 0.9958 | 0.8929 |
| C-A-L | 0.2536 | 0.9503 | 0.8095 | 0.1945 | 0.9719 | 0.8214 | 0.0739 | 0.9972 | 0.9167 |
| SG-C-A-L | 0.2314 | 0.9565 | 0.8214 | 0.1575 | 0.9823 | 0.8809 | 0.0634 | 0.9978 | 0.8809 |
| **Riverview** | | | | | | | | | |
| ARIMA | 0.4375 | 0.8132 | 0.8106 | 0.1708 | 0.9474 | 0.8038 | 0.1137 | 0.9570 | 0.7973 |
| LSTM | 0.3212 | 0.8510 | 0.8108 | 0.1537 | 0.9400 | 0.8108 | 0.0982 | 0.9705 | 0.8273 |
| ARIMA-LSTM | 0.2874 | 0.8767 | 0.8378 | 0.1314 | 0.9706 | 0.9595 | 0.0558 | 0.9934 | 0.9189 |
| SG-ARIMA-LSTM | 0.2262 | 0.9392 | 0.8243 | 0.1051 | 0.9763 | 0.8243 | 0.05639 | 0.9904 | 0.8108 |
| C-A-L | 0.2597 | 0.9483 | 0.8738 | 0.1157 | 0.9751 | 0.9324 | 0.05674 | 0.9933 | 0.9459 |
| SG-C-A-L | 0.2165 | 0.9602 | 0.8919 | 0.09214 | 0.9846 | 0.9324 | 0.05664 | 0.9939 | 0.9189 |

Note: C-A-L = CEEMDAN-ARIMA-LSTM.

**Figure 16.** The time series of observations and hybrid forecasting models for SPI prediction (Left) and their Taylor diagram plots at different timescales (Right) for SPI-6, SPI-9, and SPI-12 of False Bay Park meteorological station.

According to the statistical metrics in Table 7 and the Taylor diagram (see Figs. 11–16), the effectiveness of hybrid models that incorporate filter and signal decomposition techniques (SG and CEEMDAN) in improving prediction accuracy, particularly for drought forecasting, is highlighted. These findings support other research (Taylan et al., 2021; Elbeltagi et al., 2023; Rezaiy and Shabri, 2024), which highlights the superior accuracy of hybrid drought forecasting models compared to individual models. For example, Taylan et al. (2021) developed a hybrid model to forecast drought using precipitation data from Çanakkale, Gökçeada, and Bozcaada stations between 1975 and 2010. The study found that the hybrid models, which incorporated preprocessing techniques, performed better. Elbeltagi et al. (2023) utilized a hybrid model to estimate the SPI for 3, 6, and 12-month drought periods from 2000 to 2019. The findings demonstrated that RSS-M5P model yielded the most precise SPI predictions, with MAE = 0.497, RMSE = 0.682, RAE = 81.88, RRSE = 87.22, and $R^2 = 0.507$ for SPI-3; MAE = 0.452, RMSE = 0.717, RAE = 69.76, RRSE = 85.24, and $R^2 = 0.402$ for SPI-6 and MAE = 0.294, RMSE = 0.377, RAE = 55.79, RRSE = 59.57, and $R^2 = 0.783$ for SPI-12. The models employed to analyse drought in Jaisalmer, Rajasthan, yielded

the most effective results, exceeding those of RSS-RF and RSS-RT. Additionally, Rezaiy and Shabri (2024b) introduced a W-EEMD-ARIMA model for drought prediction. This model utilises monthly precipitation data from Kabul spanning 1970 to 2019. The $R^2$ value was 0.9946, the MAPE was 18.9674, the RMSE was 0.0736, the MAE was 0.0575, and the SPI-12 validation indicated that our model was accurate. The outcomes obtained here surpassed those of the ARIMA, Wavelet-ARIMA, and EEMD-ARIMA models in terms of raw data (RMSE: 0.0858, MAE: 0.0660, MAPE: 24.5411, $R^2$: 0.9925), analytical method (MAE: 0.1874, MAPE: 60.0220, $R^2$: 0.9361), and maximum likelihood estimation (RMSE: 0.1002, MAE: 0.0691, MAPE: 23.7122, $R^2$: 0.9898). During the SPI-3, SPI-6, and SPI-9 periods, our hybrid model consistently outperformed other models. Our proposed hybrid model surpasses ARIMA, Wavelet-ARIMA, and EEMD-ARIMA in enhancing the precision of drought predictions, as evidenced by this data.

In terms of term forecasting accuracy, the hybrid models, particularly SG-CEEMDAN-ARIMA-LSTM, consistently outperformed all other models across all SPI timescales, according to a comparison of this study's results with previous research. All models successfully reproduced the original

SPI time series. With the range values of RMSE of 0.1453–0.2314 for SPI-6, 0.0921–0.1631 for SPI-9, and 0.0349–0.07416 for SPI-12, and the highest $R^2$ values of 0.9565–0.9839 for SPI-6, 0.9836–0.9892 for SPI-9, and 0.9939–0.9990 for SPI-12 across all timescales, the SG-CEEMDAN-ARIMA-LSTM model showed the most proficiency in capturing extreme values and rapid transitions. That these methods, when combined, improve the models' capacity to represent drought in the uMkhanyakude district, both in the short and long term, is supported by the data. This makes the models far better at foretelling when droughts will occur. In light of the foregoing, our study provides useful information regarding the use of the hybrid SG-CEEMDAN-ARIMA-LSTM model to the forecasting of meteorological droughts.

## 5   Conclusions

This study examined the trends in the Standardised Precipitation Index (SPI) over different timescales (SPI-6, SPI-9, and SPI-12) utilising the Mann-Kendall (MK), modified Mann-Kendall (MMK) test, and the innovative trend analysis (ITA) protocol. The monthly rainfall data from the uMkhanyakude district, South Africa, covering the years 1980 to 2023, was used for these calculations. Rainfall has been trending downward at a 95 % confidence level, according to the MK and MMK tests. The ITA results supported these findings as well, revealing a declining trend with most of the data points going below the 1 : 1 line. To predict SPI data over various timescales, this research employed LSTM and autoregressive integrated moving average (ARIMA) models. Researchers used a hybrid model that combines the SG-CEEMDAN processing method with the ARIMA-LSTM model to enhance the precision of SPI forecasts. They also used SG filtering and full ensemble empirical mode decomposition with adaptive noise (CEEMDAN). Figures 11–16 and Table 4 display the results of a thorough comparison examination of the forecast outcomes. The results revealed that the inclusion of preprocessing techniques (SG filtering, CEEMDAN, and SG-CEEMDAN) significantly improved the model performance in forecasting SPI at all timescales. The performance consistently increased with higher timescales, potentially due to lower noise levels. Across different timescales, the SG and CEEMDAN combined hybrid model consistently outperformed the individual models. Notably, the CEEMDAN-ARIMA-LSTM model outperformed the SG-ARIMA-LSTM model at all timescales, while the SG-CEEMDAN-ARIMA-LSTM model consistently exhibited the lowest root mean square error (RMSE) values across all indices. These results demonstrate that combining SG-CEEMDAN with ARIMA-LSTM has the potential to significantly enhance the accuracy of meteorological drought forecasting.

The principal conclusion of the study is that ARIMA-LSTM, in conjunction with SG, CEEMDAN, and SG-CEEMDAN, serves as an effective instrument for early warning systems and meteorological drought prediction. The proposed methodology in this paper serves as a framework for modeling complex meteorological phenomena such as drought, which is particularly pertinent in semi-arid regions. Enhancing model performance and creating efficient models for weather forecasting can be achieved through techniques that address data noise, nonlinearity, and nonstationarity. To enhance water resource management, make informed decisions regarding agricultural output and tourism management, and establish regulations, it is essential to acquire extremely effective models for drought prediction. The omission of exogenous environmental variables in the SG-CEEMDAN-ARIMA-LSTM model represents a significant drawback of the study. The model's forecast accuracy and real-world application are limited by disregarding these exogenous effects, which can substantially affect drought conditions. Future studies should aim to include external variables, including temperature, soil moisture, vegetation indices, and anthropogenic factors such as land use and water management, to improve the model's efficacy. This integration would provide a more thorough comprehension of drought dynamics, hence improving the model's accuracy and dependability in drought predictions. Additionally, it is essential to investigate alternate decomposition methods, such as enhanced CEEMDAN (iCEEMDAN), which may provide significant insights.

# References

Alashan, S.: An improved version of innovative trend analyses, Arab. J. Geosci., 11, 50, https://doi.org/10.1007/s12517-018-3393-x, 2018.

Alashan, S.: Combination of modified Mann-Kendall method and Şen innovative trend analysis, Eng. Rep., 2, e12131, https://doi.org/10.1002/eng2.12131, 2020.

Alquraish, M., Abuhasel, K. A., Alqahtani, S. A., and Khadr, M.: SPI-based hybrid hidden Markov–GA, ARIMA–GA, and ARIMA–GA–ANN models for meteorological drought forecasting, Sustainability, 13, 12576, https://doi.org/10.3390/su132212576, 2021.

Ashraf, M. S., Shahid, M., Waseem, M., Azam, M., and Rahman, K. U.: Assessment of variability in hydrological droughts using the improved innovative trend analysis method, Sustain., 15, 9065, https://doi.org/10.3390/su15119065, 2023.

Bagmar, M. S. H. and Khudri, M. M.: Application of box-jenkins models for forecasting drought in north-western part of Bangladesh, Environmental Engineering Research, 26, https://doi.org/10.4491/eer.2020.294, 2021.

Balti, H., Abbes, A. B., Mellouli, N., Farah, I. R., Sang, Y., and Lamolle, M.: A review of drought monitoring with big data: Issues, methods, challenges, and research directions, Ecological Informatics, 60, 101136, https://doi.org/10.1016/j.ecoinf.2020.101136, 2020.

Bard, A., Renard, B., Lang, M., Giuntoli, I., Korck, J., Koboltschnig, G., Janža, M., d'Amico, M., and Volken, D.: Trends in the hydrologic regime of Alpine rivers, J. Hydrol., 529, 1823–1837, https://doi.org/10.1016/j.jhydrol.2015.08.052, 2015.

Box, G. E., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M.: Time series analysis: forecasting and control, 5th Edn., John Wiley & Sons, Hoboken, NJ, https://doi.org/10.1111/jtsa.12194, 2015.

Bukhosini, Z. and Moyo, I.: An analysis of the challenges faced by small-scale farmers and their response to the 2014–2016 drought in Mfekayi, Mtubatuba, KZN, South Africa, Afr. J. Dev. Stud., 13, 1, 2023.

Caloiero, T., Coscarelli, R., Ferrari, E., and Mancini, M.: Trend detection of annual and seasonal rainfall in Calabria (Southern Italy), Int. J. Climatol., 31, 44–56, https://doi.org/10.1002/joc.2054, 2011.

Ding, Y., Yu, G., Tian, R., and Sun, Y.: Application of a hybrid CEEMD-LSTM model based on the standardized precipitation index for drought forecasting: the case of the Xinjiang Uygur Autonomous Region, China, Atmosphere, 13, 1504, https://doi.org/10.3390/atmos13091504, 2022.

Elbeltagi, A., Kumar, M., Kushwaha, N. L., Pande, C. B., Ditthakit, P., Vishwakarma, D. K., and Subeesh, A.: Drought indicator analysis and forecasting using data driven models: case study in Jaisalmer, India, Stoch. Environ. Res. Risk Assess., 37, 113–131, https://doi.org/10.1007/s00477-022-02277-0, 2023.

Gudko, V., Tanwar, S., Minkina, T., Sushkova, S., Usatov, A., Azarin, K., Safronenkova, I., Melnik, Y., Voloshchuk, V., Gülser, C., and Kızılkaya, R.: Analysis of drought dynamics using SPI and SARIMA models: A case study of the Rostov Region, Russia, Eurasian Journal of Soil Science 14, 208–218, https://doi.org/10.18393/ejss.1682888, 2025.

Hamed, K. H. and Rao, A. R.: A modified Mann-Kendall trend test for autocorrelated data, J. Hydrol., 204, 182–196, https://doi.org/10.1016/S0022-1694(97)00125-X, 1998.

Harka, A. E., Jilo, N. B., and Behulu, F.: Spatial-temporal rainfall trend and variability assessment in the Upper Wabe Shebelle River Basin, Ethiopia: Application of innovative trend analysis method, J. Hydrol. Reg. Stud., 37, 100915, https://doi.org/10.1016/j.ejrh.2021.100915, 2021.

Helsel, D. R. and Hirsch, R. M.: Statistical methods in water resources, Elsevier, Amsterdam, https://doi.org/10.3133/tm4A3, 1993.

Hırca, T., Eryılmaz Türkkan, G., and Niazkar, M.: Applications of innovative polygonal trend analyses to precipitation series of Eastern Black Sea Basin, Turkey, Theor. Appl. Climatol., 147, 651–667, https://doi.org/10.1007/s00704-021-03837-0, 2022.

Hochreiter, S. and Schmidhuber, J.: Long short-term memory, Neural Comput., 9, 1735–1780, https://doi.org/10.1162/neco.1997.9.8.1735, 1997.

Hussain, A., Rizwan, N., Al-Rezami, A. Y., Adam M. O., Fuad S. A., and Mohammed M. A. A.: Application of random forest for identification of an appropriate model for predicting meteorological drought, Advances in Meteorology, 2025, 7674140, https://doi.org/10.1155/adme/7674140, 2025.

Kalisa, W., Zhang, J., Igbawua, T., Kayiranga, A., Ujoh, F., Aondoakaa, I. S., and Nibagwire, D.: Spatial multi-criterion decision making (SMDM) drought assessment and sustainability over East Africa from 1982 to 2015, Remote Sensing, 13, 5067, https://doi.org/10.3390/rs13245067, 2021.

Kendall, M.: Rank Correlation Methods, Griffin, London, https://www.cabidigitallibrary.org/doi/full/10.5555/19521603271, 1975.

Kganyago, M., Mukhawana, M. B., Mashalane, M., Mgabisa, A., and Moloele, S.: Recent trends of drought using remotely sensed and in-situ indices: Towards an integrated drought monitoring system for South Africa, in: 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, 6225–6228, https://doi.org/10.1109/IGARSS47720.2021.9553994, 2021.

Latifoğlu, L. and Özger, M.: A novel approach for high-performance estimation of SPI data in drought prediction, Sustainability, 15, 14046, https://doi.org/10.3390/su151914046, 2023.

Likinaw, A., Alemayehu, A., and Bewket, W.: Trends in extreme precipitation indices in Northwest Ethiopia: Comparative analysis using the Mann–Kendall and innovative trend analysis methods, Climate, 11, 164, https://doi.org/10.3390/cli11080164, 2023.

Lloyd-Hughes, B. and Saunders, M. A.: A drought climatology for Europe, Int. J. Climatol., 22, 1571–1592, https://doi.org/10.1002/joc.846, 2002.

Ma, X., He, Y., Xu, J., van Noordwijk, M., and Lu, X.: Spatial and temporal variation in rainfall erosivity in a Himalayan watershed, Catena, 121, 248–259, https://doi.org/10.1016/j.catena.2014.05.012, 2014.

Malik, A., Kumar, A., Rai, P., and Kuriqi, A.: Prediction of multi-scalar standardized precipitation index by using artificial intelligence and regression models, Climate 9, 28, https://doi.org/10.3390/cli9020028, 2021.

Mann, H. B.: Nonparametric tests against trend, Econometrica, 13, 245–259, https://doi.org/10.2307/1907187, 1945.

McKee, T. B., Doesken, N. J., and Kleist, J.: The relationship of drought frequency and duration to time scales, in: Proceedings

https://doi.org/10.5194/nhess-26-315-2026

Nat. Hazards Earth Syst. Sci., 26, 315–342, 2026

of the 8th Conference on Applied Climatology, Anaheim, California, American Meteorological Society, 17, 179–183, 1993.

McKee, T. B., Doesken, N. J., and Kleist, J.: Drought monitoring with multiple time scales, in: Proceedings of the Conference on Applied Climatology, Boston, MA, USA, American Meteorological Society, 1995.

Mirabbasi, R., Ahmadi, F., and Jhajharia, D.: Comparison of parametric and non-parametric methods for trend identification in groundwater levels in Sirjan plain aquifer, Iran, Hydrol. Res., 51, 1455–1477, https://doi.org/10.2166/nh.2020.041, 2020.

Montgomery, D. C., Jennings, C. L., and Kulahci, M.: Introduction to time series analysis and forecasting, 2nd Edn., Wiley, Hoboken, NJ, 2015.

Naik, M. and Abiodun, B. J.: Projected changes in drought characteristics over the Western Cape, South Africa, Meteorological Applications, 27, e1802, https://doi.org/10.1002/met.1802, 2020.

Ngwenya, M., Gidey, E., and Simatele, M. D.: Agroecological-based modeling of meteorological drought at 12-month time scale in the Western Cape Province of South Africa, Earth Sci. Inform., 17, 1851–1865, https://doi.org/10.1007/s12145-023-01193-3, 2024.

Öztopal, A. and Şen, Z.: Innovative trend methodology applications to precipitation records in Turkey, Water Resour. Manag., 31, 727–737, https://doi.org/10.1007/s11269-016-1343-5, 2017.

Rezaiy, R. and Shabri, A.: An innovative hybrid W-EEMD-ARIMA model for drought forecasting using the standardized precipitation index, Nat. Hazards, https://doi.org/10.1007/s11069-024-06758-z, 2024b.

Savitzky, A. and Golay, M. J. E.: Smoothing and differentiation of data by simplified least squares procedures, Anal. Chem., 36, 1627–1639, https://doi.org/10.1021/ac60214a047, 1964.

Şen, Z.: Innovative trend analysis methodology, J. Hydrol. Eng., 17, 1042–1046, https://doi.org/10.1061/(ASCE)HE.1943-5584.0000556, 2012.

Şen, Z.: Innovative trend significance test and applications, Theor. Appl. Climatol., 127, 939–947, https://doi.org/10.1007/s00704-015-1681-x, 2017.

Sharifi, A., Baubekova, A., Patro, E. R., Klöve, B., and Haghighi, A. T.: The combined effects of anthropogenic and climate change on river flow alterations in the Southern Caspian Sea, Iran, Heliyon, 10, e18663, https://doi.org/10.1016/j.heliyon.2024.e31960, 2024.

Sibiya, S., Mbatha, N., Ramroop, S., Melesse, S., and Silwimba, F.: Forecasting of Standardized Precipitation Index Using Hybrid Models: A Case Study of Cape Town, South Africa, Water, 16, 2469, https://doi.org/10.3390/w16172469, 2024.

Song, Y. and Park, M.: A study on the appropriateness of the drought index estimation method using damage data from Gyeongsangnamdo, South Korea, Atmosphere, 12, 998, https://doi.org/10.3390/atmos12080998, 2021.

Tan, Y. X., Ng, J. L., and Huang, Y. F.: A review on drought index forecasting and their modelling approaches, Archives of Computational Methods in Engineering, 30, 1111–1129, https://doi.org/10.1007/s11831-022-09828-2, 2023.

Taylan, E. D.: An approach for future droughts in Northwest Türkiye: SPI and LSTM methods, Sustainability 16, 6905, https://doi.org/10.3390/su16166905, 2024.

Taylan, E. D., Özlem, T., and Baykal, T.: Hybrid wavelet-artificial intelligence models in meteorological drought estimation, Journal of Earth System Science 130, 38, https://doi.org/10.1007/s12040-020-01488-9, 2021.

Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, J. Geophys. Res.-Atmos., 106, 7183–7192, https://doi.org/10.1029/2000JD900719, 2001.

Wang, X., Hou, X., and Wang, Y.: Spatiotemporal variations and regional differences of extreme precipitation events in the coastal area of China from 1961 to 2014, Atmos. Res., 197, 94–104, https://doi.org/10.1016/j.atmosres.2017.06.010, 2017.

Wilhite, D. A. and Glantz, M. H.: Understanding: the drought phenomenon: the role of definitions, Water International, 10, 111–120, https://doi.org/10.1080/02508068508686328, 1985.

Xu, D., Ding, Y., Liu, H., Zhang, Q., and Zhang, D.: Applicability of a CEEMD–ARIMA combined model for drought forecasting: a case study in the Ningxia Hui Autonomous Region, Atmosphere, 13, 1109, https://doi.org/10.3390/atmos13071109, 2022.

Yue, S., Pilon, P., Phinney, B., and Cavadias, G.: The influence of autocorrelation on the ability to detect trend in hydrological series, Hydrol. Process., 16, 1807–1829, https://doi.org/10.1002/hyp.1095, 2002.

Zena, B. K., Demissie, T. A., and Feyessa, F. F.: Comparative analysis of long-term precipitation trends and its implication in the Modjo catchment, central Ethiopia, J. Water Clim. Change, 13, 3883–3905, https://doi.org/10.2166/wcc.2022.234, 2022.

Zhang, G. P.: Time series forecasting using a hybrid ARIMA and neural network model, Neurocomputing, 50, 159–175, https://doi.org/10.1016/S0925-2312(01)00702-0, 2003.

Zhang, H., Loaiciga, H. A., and Sauter, T.: A novel fusion-based methodology for drought forecasting, Remote Sensing, 16, 828, https://doi.org/10.3390/rs16050828, 2024.

Zhang, X., Duan, Y., Duan, J., Chen, L., Jian, D., Lv, M., and Ma, Z.: A daily drought index-based regional drought forecasting using the Global Forecast System model outputs over China, Atmospheric Research, 273, 106166, https://doi.org/10.1016/j.atmosres.2022.106166, 2022.

Zhang, X., Qiao, W., Huang, J., Shi, J., and Zhang, M.: Flow prediction in the lower Yellow River based on CEEMDAN-BILSTM coupled model, Water Supply, 23, 396–409, https://doi.org/10.2166/ws.2022.215, 2023.

Zuo, D., Hou, W., Wu, H., Yan, P., and Zhang, Q.: Feasibility of calculating standardized precipitation index with short-term precipitation data in China, Atmosphere, 12, 603, https://doi.org/10.3390/atmos12050603, 2021.