



Contents lists available at ScienceDirect

Engineering Applications of Artificial Intelligence

journal homepage: www.elsevier.com/locate/engappai

Research paper

Establishing hybrid deep learning models for regional daily rainfall time series forecasting in the United Kingdom

Geethu Thottungal Harilal^a, Aniket Dixit^{a,b}, Giovanni Quattrone^{a,c,*}^a Middlesex University, Department of Computer Science, UK^b Coventry University, Computational Science and Mathematical Modelling Research Centre, UK^c University of Turin, Department of Computer Science, Italy

ARTICLE INFO

Keywords:

Deep learning
Long short term memory
Recurrent neural networks
Convolutional neural networks
Daily rainfall forecasting

ABSTRACT

Accurate daily rainfall predictions are becoming increasingly important, particularly in the era of changing climate conditions. These predictions are essential for various sectors, including agriculture, water resource management, flood preparedness, and pollution monitoring. This study delves into the complex relationship between meteorological data, with a focus on the accurate forecasting of rainfall by identifying the impact of temperature variations on rainfall patterns in different regions of the United Kingdom (UK). The meteorological data was collected from the National Aeronautics and Space Administration (NASA) and covers daily observations from January 1, 1981, to July 31, 2023, in four distinct regions of the UK: England, Wales, Scotland, and Northern Ireland. The main objective of this research is to introduce hybrid deep learning models, namely Convolutional Neural Networks (CNN) with Long Short Term Memory (LSTM) and Recurrent Neural Networks (RNN) with Long Short Term Memory (LSTM), for predicting daily rainfall using time-series data from the four UK countries, specifically designed for daily rainfall forecasting of four regions in the UK. The models are fine-tuned using the hyperparameter optimisation method. Comprehensive performance evaluations, including Loss Function, Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE), are employed to compare the effectiveness of our proposed hybrid models with established baseline models, including LSTM, stacked LSTM, and Bidirectional LSTM. Additionally, a visual analysis of actual and predicted rainfall data is conducted to identify the most proficient forecasting model for each region. Results reveal that the proposed hybrid models consistently outperform other models in terms of both quantitative performance metrics and visual assessments across all four regions in the UK. This research contributes to improved rainfall forecasting methodologies, which are critical for sustainable agricultural practices and resource management.

1. Introduction

Rainfall is a pivotal meteorological factor, and the UK is renowned for its frequent and often excessive precipitation. This weather phenomenon significantly impacts the UK's economy and agriculture, as it can lead to devastating floods or crippling droughts (Trenberth, 2005). Wetter regions in the UK experience over 200 rainy days per year, encompassing more than half of the calendar, while drier areas witness an annual average of 150 to 200 rainy days (GOV.UK, 2023). Over the past decade, the UK has seen an increase in the number of heavy rain days, surpassing the thresholds set at 95% and 99% of the average rainfall from 1961 to 1990. Furthermore, occurrences of rainfall exceeding 50 mm have become more frequent, signalling intensification and increased frequency of precipitation across the UK (Cotterill et al., 2021).

The impact of human-induced greenhouse gas emissions on global temperatures is widely acknowledged, resulting in significant temperature increases in the UK over the past few decades, particularly in the last 20 years (GOV.UK, 2023). As human activities contribute to climate change, it becomes imperative to investigate the relationship between temperature changes and rainfall in different regions in the UK for efficient adaptation and disaster mitigation strategies (Nakicenovic et al., 2000). Rainfall is a complex process influenced by various meteorological factors and conditions over different time scales, ranging from one day to several days (Barrera-Animas et al., 2022). A comprehensive understanding of these factors is critical for managing and predicting rainfall accurately.

The prediction of rainfall is a crucial element of weather forecasting. In earlier studies, this has been accomplished by using statistical

* Corresponding author.

E-mail addresses: gt382@live.mdx.ac.uk (G. Thottungal Harilal), dixita4@uni.coventry.ac.uk (A. Dixit), g.quattrone@mdx.ac.uk (G. Quattrone).

<https://doi.org/10.1016/j.engappai.2024.108581>

Received 15 December 2023; Received in revised form 29 February 2024; Accepted 4 May 2024

Available online 21 May 2024

0952-1976/© 2024 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY-NC license (<http://creativecommons.org/licenses/by-nc/4.0/>).

techniques that attempt to correlate rainfall with geographic coordinates and various atmospheric factors. However, these methods often struggle to accurately predict rainfall patterns due to their inherent complexity and nonlinearity (Wu and Chau, 2013). In recent years, advanced techniques such as Singular Spectrum Analysis, Empirical Mode Decomposition, and Wavelet analysis have been explored to address these challenges (Gan et al., 2018). However, these approaches can be computationally intensive and may only yield limited improvements in forecast accuracy (Singh and Borah, 2013).

Previous studies conducted to predict rainfall in the UK have predominantly concentrated on specific forecasts and have not sufficiently delved into the regional diversity of the country (Barrera-Animas et al., 2022). This drawback underscores the necessity for a more exhaustive methodology that takes into account the distinct geographical and meteorological attributes present in different regions across the UK. The emergence of Artificial Neural Networks (ANNs) has revolutionised rainfall forecasting by offering a more flexible and adaptable approach. Among various types of ANNs, Recurrent Neural Networks (RNNs) have become popular due to their ability to address temporal dynamics that are present in meteorological time-series data (Liu et al., 2019). However, traditional RNNs have limitations in learning and making accurate long-term forecasts (Ni et al., 2020). Variants of RNNs, such as Long Short-Term Memory (LSTM) Networks, have been developed to overcome their limitations. LSTM Networks are equipped with memory cells that retain information over extended time periods and have shown better performance in multi-step ahead predictions compared to traditional RNNs, as demonstrated by various studies (Greff et al., 2016; Kratzert et al., 2018; Yunpeng et al., 2017).

Rainfall forecasting using ANNs has promising potential, but challenges still exist, especially in capturing regional rainfall patterns' spatial and temporal variations (Hossain et al., 2020). To enhance forecasting accuracy, advanced architectures such as Bidirectional LSTM Networks have been proposed, which leverage information from both past and future sequential data (Balluff et al., 2020). Design decisions and model implementation play a critical role in determining the effectiveness of ANN-based forecasting models, despite advancements in this field (Hutter et al., 2019). This study aims to bridge existing gaps in rainfall forecasting methods, ultimately improving accuracy and enhancing understanding of the impact of temperature variations on rainfall patterns across the UK. In this study we provide the following main contributions:

- Investigate the impact of various meteorological parameters, with a focus on temperature, on rainfall, and explore trends and antecedent effects.
- Adapt three LSTM-based models, used as benchmarks – namely LSTM, Stacked-LSTM, and Bidirectional LSTM Networks – using time-series data from the four UK countries to predict daily rainfall.
- Propose two hybrid models, CNN with LSTM (CLSTM) and RNN with LSTM (RLSTM), for predicting daily rainfall using time-series data from the four UK countries.
- Evaluate the performance of each model in terms of their ability to forecast daily rainfall amounts using time-series data from the four UK countries.

This study reveals that both CLSTM and RLSTM hybrid models offer consistent superiority over LSTM, Stacked LSTM, and Bidirectional LSTM in UK rainfall prediction across all UK regions. Specifically, RLSTM adeptly captures sequential dependencies and long-term patterns, offering reliable forecasts adaptable to diverse weather dynamics. In contrast, LSTM, stacked LSTM, and Bidirectional LSTM encounter limitations in handling intricate temporal patterns. RLSTM emerges as the top choice due to its robustness and adaptability.

The rest of this paper is structured as follows: First, the detailed background and previous studies are discussed in Section 2. The Section 3 provides an in-depth exploration of Neural Networks, Recurrent

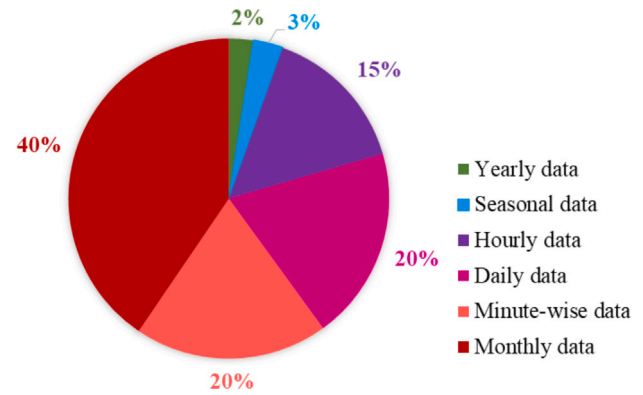


Fig. 1. Percentage of datasets categorised by their projected time frame (Hussein et al., 2022).

Neural Networks, including CNN, RNN, LSTM, and Bidirectional LSTM, which form the basis of the models in this research. Section 4 details the study area, dataset description, and methodology employed for exploratory data analysis, feature selection, and the development of the proposed hybrid deep learning models. This includes the steps for pre-processing data, the rationale for proposed architectures, and hyperparameter tuning. Subsequently, Section 5 is devoted to presenting and discussing the obtained results. Finally, Section 6 encapsulates the conclusions drawn from the findings and outlines future research endeavours.

2. Related work

Over the last two decades, an extensive body of research has been dedicated to the prediction of rainfall time series utilising machine learning and deep learning techniques. During the period from 1996 to 2014, researchers frequently employed various iterations of ANNs to enhance prediction accuracy. These ANN models encompassed a range of input parameters, including Min–Max temperature, relative humidity, average humidity, average wind speed, wind direction, latitude–longitude coordinates, sea surface pressure, and rainfall (Parmar et al., 2017). Notably, among these parameters, Min–Max temperature and humidity played a prominent role in rainfall prediction modelling.

Between 2016 and 2020, deep learning models such as CNN, LSTM, Convolutional LSTM (ConvLSTM), and RNN garnered substantial attention in the field of rainfall time series forecasting due to their heightened predictive accuracy. Temperature, humidity, wind speed, and air pressure emerged as the most frequently considered parameters for forecasting purposes (Hussein et al., 2022). The short-term forecasting of precipitation involves predicting future rainfall within a relatively brief timeframe, ranging from a few hours to several days. ConvLSTM, renowned for its ability to integrate spatio-temporal information, proves highly suitable for addressing such forecasting challenges (Wu et al., 2022). Conversely, long-term forecasting extends beyond the short-range, encompassing periods of weeks to months or even years (Markuna et al., 2023). Fig. 1 illustrates that a majority of the datasets assessed in prior research predominantly exhibit monthly temporal granularity. Input features employed in these investigations can be broadly categorised into two groups: 1D input features, where each time lag corresponds to one or more geophysical parameters recorded at established, fixed locations such as weather stations, and 2D input features, where each time lag encapsulates a 2D spatial representation of precipitation values across the geographical study area, often collected through satellite imaging.

The process employed in earlier studies focused on rainfall prediction involves several key stages. The initial step involves collecting meteorological data from various monitoring sites situated within the

target region. The raw data is subsequently subjected to rigorous cleansing and filtration procedures to ensure its quality and reliability. This preparatory phase plays a pivotal role in laying the foundation for subsequent analysis. Following data cleansing, researchers typically proceed to the feature selection stage, employing techniques such as correlation analysis to identify the most relevant meteorological parameters and variables. The selection of pertinent features is a crucial decision, as it directly influences the accuracy and performance of forecasting models. Subsequent to feature selection, the chosen features are incorporated into machine learning techniques for rainfall prediction. Deep learning models, particularly LSTM and its derivatives, are frequently employed for short-term rainfall forecasting tasks, owing to their proficiency in capturing temporal patterns and dependencies, rendering them well-suited for short-range predictions. In the context of longer-term forecasting, researchers often opt for alternative machine learning techniques, including ANN, Support Vector Regression (SVR), and Random Forest (RF). These models offer robust capabilities for modelling complex, protracted rainfall trends and patterns. It is noteworthy that hybrid models, such as Auto-Regressive Integrated Moving Average (ARIMA), RNN, and CNN, have also garnered attention in the literature, harnessing the strengths of diverse modelling techniques to enhance the accuracy and versatility of rainfall forecasts (Hussein et al., 2022).

Kim and Bae (2017) proposed a model based on LSTM for the Gangneung region in Korea, utilising weather data from 2012, which included temperature, wind speed, humidity, and sea surface pressure. This model also took into consideration the lag characteristics of past and current hours' observations of rainfall amounts as features. Comparative analysis between the LSTM-Networks model and the ANN model revealed superior performance in terms of the RMSE evaluation metric for the proposed LSTM model. In a study conducted by Chao et al. (2018), five models were juxtaposed for predicting rainfall amounts in the Wuhan region of China, comprising ARIMA, RF, Backpropagation Neural Networks, Support Vector Machine (SVM), and LSTM-Networks. These models were evaluated using weather data from 2015 and 2016, encompassing features such as wind speed, wind direction, temperature, humidity, pressure, rainfall amount, and radiation. The LSTM-Networks model outperformed the other models in terms of both RMSE and MAE metrics. Kumar et al. (2019) conducted a study that compared two models utilising RNNs and LSTM-Networks for forecasting monthly rainfall in India from 1871 to 2016. They employed a climate dataset containing the average monthly rainfall, incorporating lag features of rainfall for the preceding 12 months. Through comprehensive evaluation, it was established that the LSTM-Networks model exhibited superior performance across multiple evaluation metrics, including RMSE, correlation coefficient (R), Nash-Sutcliffe Efficiency (NSE), and MAE.

Historically, previous studies have encountered challenges in accurately predicting precipitation, leading to issues such as overfitting, inaccurate predictions on test and validation datasets, and an inability to capture peak values. Consequently, this research endeavours to fine-tune prediction models to narrow the disparity between predicted and actual values. It is imperative to discern the relationship between temperature variables and rainfall during the data analysis phase to identify pertinent features effectively. Meteorological datasets typically furnish temperature values, encompassing minimum, maximum, and average temperatures, albeit prior studies have variably adopted distinct temperature values. Some have utilised maximum temperature values (Venkata Ramana et al., 2013; Haidar and Verma, 2016), while others have relied on average values (Hernández et al., 2016). Some studies have amalgamated minimum and maximum temperatures (Ramsundram et al., 2016; He et al., 2022), while others have incorporated all three (Saikhu et al., 2017; Xu et al., 2020). Nevertheless, conclusive evidence is lacking regarding the temperature value possessing the utmost significance in forecasting. To ascertain the correct temperature value and its antecedent impact on rainfall prediction,

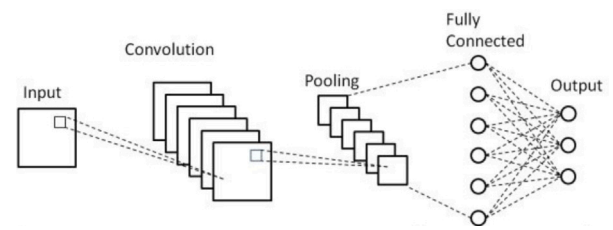


Fig. 2. CNN architecture.

a region-specific analysis will be undertaken in this proposed study. The central aim of this research is to conduct an exhaustive analysis aimed at identifying the most appropriate temperature features and their antecedent effects, with the overarching objective of enhancing rainfall forecasting precision.

3. Deep learning models

This section provides an overview of essential neural network concepts, specifically focusing on RNNs, LSTM networks, Stacked-LSTM networks, and Bidirectional LSTM networks. Previous studies have demonstrated the efficacy of neural networks, particularly LSTM-based networks, in rainfall forecasting tasks (Barrera-Animas et al., 2022; Chao et al., 2018; Kim and Bae, 2017; Kratzert et al., 2018; Kumar et al., 2019; Poornima and Pushpalatha, 2019). These neural networks provide a robust framework for weather forecasting due to their ability to handle uncertainty, capture spatiotemporal dependencies, and efficiently model discrete distributions.

3.1. CNN

CNNs were developed for processing and analysing data in a grid-like format, particularly suited for tasks involving images or sequences of data (Krizhevsky et al., 2012; Wang et al., 2017). CNNs have had a profound impact on the field of computer vision, demonstrating exceptional performance in tasks such as image classification, object detection, and image segmentation. Fig. 2 illustrates a typical deep CNN architecture, which consists of essential components such as convolutional layers, pooling layers, and a fully connected layer. Among these components, convolutional layers are of greater importance (Haidar and Verma, 2016). These convolutional layers enable CNNs to learn hierarchical features from data by connecting neurons with local regions in preceding layers, rather than all neighbouring neurons, rendering CNNs highly effective for processing visual information as well as time series data. CNN effectively captures local patterns and spatial dependencies in time series data. It can automatically learn hierarchical features from raw data and is robust to noise while being invariant to translations in the time domain.

CNNs use a filter (or kernel) matrix to analyse images or sequences of data. This matrix slides through blocks of the input layer, forming the convoluted layer. The resulting pixel of the convoluted layer is calculated using the following equation (Ojo et al., 2019):

$$C_k = f(x * W + b) \quad (1)$$

where C_k is the k th pixel of the convoluted layer, x is the corresponding pixel value, W is the coefficient vector, $f(\cdot)$ is the activation function and b is the bias. The pooling layer is responsible for down-sampling an image. To calculate the k th value of the pooling layer, the following formula is used:

$$P_k = f(\beta * \text{down}(C) + \alpha) \quad (2)$$

In this equation P_k is the k th value in the pooling layer matrix, C is the value vector from the convoluted layer, β is the coefficient, and α is the bias.

The pooling layer is used to down-sample an image. One commonly used method for pooling is called Max Pooling, where the maximum value is calculated for a block in the matrix. The calculation for the Max Pooling value is as follows:

$$\text{down}(C) = \max \{C_{s,l} \mid |s| \leq \frac{m}{2}, |l| \leq \frac{m}{2}, s, l \in \mathbb{Z}^+\} \quad (3)$$

In this equation $C_{s,l}$ is the pixel value at C of the matrix, m is the size for sub-sampling, \mathbb{Z}^+ represents the set of positive integers.

3.2. RNN

RNNs have proven invaluable in various fields, including machine translation, sentiment analysis, speech recognition, and weather forecasting, owing to their capacity to handle sequential input. By introducing a hidden state and preserving a relationship between previous and current observations, RNN cells capture data dependencies. Eq. (4) illustrates how these dependencies are maintained in time-series sequences:

$$h_t = (h_{t-1}, X_t) \quad (4)$$

where h_t and h_{t-1} are hidden states at times t and $t - 1$, respectively; X_t is the current input value at time t .

While RNNs excel at learning sequential data, they encounter challenges in grasping long-range dependencies, resulting in vanishing error gradients during backward propagation. However, variants such as LSTM-Networks have emerged to address this issue. LSTM networks incorporate memory cells and gating mechanisms, enabling them to effectively manage long-term dependencies, making them valuable for processing sequential data (Singh et al., 2015; Barrera-Animas et al., 2022).

3.3. LSTM

Among the various RNN variants, LSTM Networks are the most popular due to their ability to capture longer dependencies within sequential data. They have been successful in diverse research domains, including weather prediction, speech recognition, traffic forecasting, and human trajectory prediction. LSTMs employ three gates – input, output, and forget gates – to learn dependencies from the most recent two states and the current state, effectively addressing the vanishing gradient problem. The input gate controls the amount of new state information used, the output gate determines information retention from earlier states, and the forget gate regulates information flow within the internal state. The core LSTM definitions are provided in Eq. (5), and the LSTM block diagram is represented in Fig. 3.

$$\begin{aligned} \text{input gate } (i) &= (W_i h_{t-1} + U_i X_t + V_i C_{t-1}) \\ \text{output gate } (o) &= (W_o h_{t-1} + U_o X_t + V_o C_{t-1}) \\ \text{forget gate } (f) &= (W_f h_{t-1} + U_f X_t + V_f C_{t-1}) \\ \text{internal hidden state } (g) &= \tanh(W_g h_{t-1} + U_g X_t) \\ \text{current cell state } c_t &= (f * c_{t-1}) + (g * i) \\ \text{current hidden state } h_t &= \tanh(c_t) * o \end{aligned} \quad (5)$$

where W, U, V represent different weight matrices, t is current time and $t - 1$ refers to the preceding time step.

The Stacked-LSTM Networks extend this architecture by sequentially connecting several LSTM Networks to achieve a deeper representation of time-series data as illustrated in Fig. 4 (Cui et al., 2018). The stacked LSTM can capture complex temporal patterns through multiple layers of recurrent units. This allows for a hierarchical feature representation, where each layer captures different levels of abstraction. The stacked LSTM also provides flexibility in model architecture, allowing for the incorporation of domain-specific knowledge.

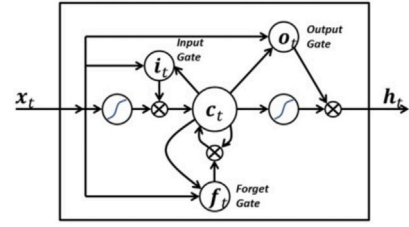


Fig. 3. LSTM architecture.

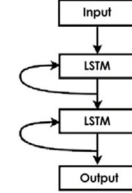


Fig. 4. Stacked LSTM architecture.

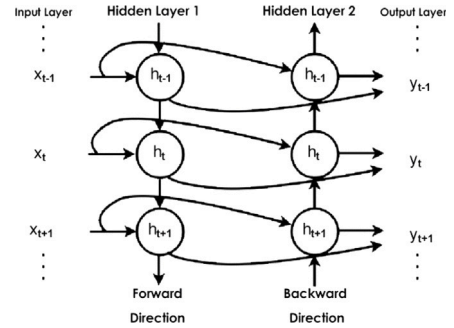


Fig. 5. Bidirectional LSTM architecture.

3.4. Bidirectional LSTM

Bidirectional RNNs are a variant of traditional RNNs that excel in capturing dependencies from both past and future states, making them particularly effective in natural language processing, speech recognition, time series forecasting, and handwriting recognition. This architecture employs two types of RNN cells: one captures data from left to right (standard RNN cells), and the other captures data in reverse. This combined approach enhances the model's ability to comprehend sequential data comprehensively (Cui et al., 2018; Cheng et al., 2019). When LSTM cells replace the RNN cells in a Bidirectional RNN, it becomes a Bidirectional LSTM Network. The structure of a Bidirectional LSTM network is depicted in Fig. 5. The Bidirectional LSTM is effective in capturing long-term dependencies. This is due to its memory cell structure which enables it to process both past and future information through forward and backward processing. It is especially suitable for sequential data such as time series where past and future context is important (Barrera-Animas et al., 2022).

4. Materials and methods

This section describes a detailed overview of the methodology (illustrated in Fig. 6) utilised to predict rainfall in four distinct regions of the UK: England, Wales, Scotland, and Northern Ireland. The process includes collecting meteorological data from four regions of UK, pre-processing, conducting exploratory data analysis, feature selection, and developing rainfall prediction models using deep learning.

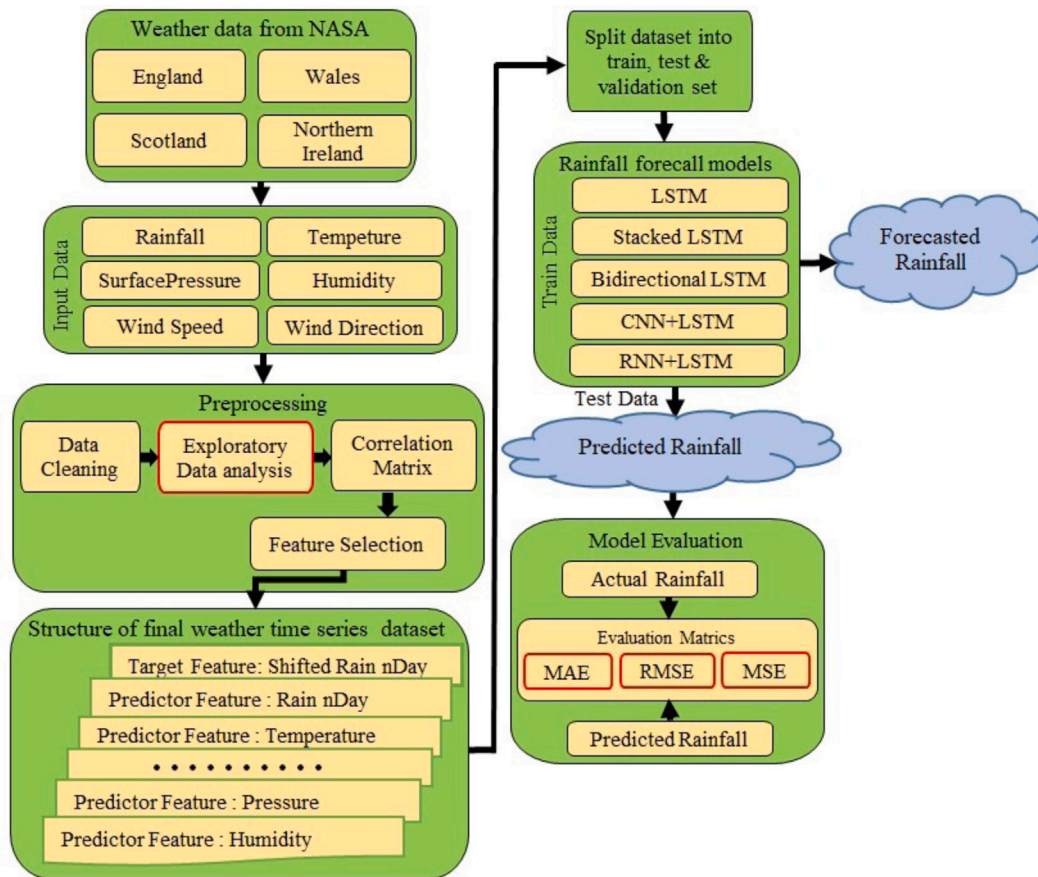


Fig. 6. Architecture of the proposed model of rainfall forecasting.

4.1. Study area and data description

The study focuses on the data from four geographical regions or constituent countries of the UK: England, Scotland, Wales, and Northern Ireland, all located in the British Isles. Northern Ireland features a mild, humid climate with mild winters and cool summers, while Scotland is known for its cold and rainy conditions, particularly in mountainous areas. England experiences an unstable climate with occasional fogs and stormy winds, and Wales shares a similar climate with variations in rainfall. Coastal areas in the UK are susceptible to sea fog due to the interaction of the Gulf Stream and cold Atlantic air, rendering the region's weather one of the most unpredictable in Europe (Neagh and Lomond, 2014). The study explores the unique climatic conditions and variations in these four UK regions.

Historical data spanning from 1st January 1981 to 31st July 2023 for the four UK regions, was obtained from NASA Power Data Viewer for this study (Power, 2023). The dataset is accessible from the link: <https://power.larc.nasa.gov/data-access-viewer/>. The dataset comprises 24 columns, including year, month, and date of the month columns, encompassing 15,551 rows of entries in four separate dataset for each region and the description of each field is given in Table 1. It encompasses daily recorded weather measurements, including rainfall, temperature, various temperature parameters, surface pressure, relative and specific humidity, as well as wind speed and direction. The temperature parameters at 2 m above ground level is a standard measure in meteorology because it reflects the near-surface air temperature that humans and ecosystems experience. This standard measurement is widely employed in weather prediction and climatological studies due to its direct influence on evaporation rates, atmospheric stability, and cloud formation processes (Quej et al., 2022; Ben Bouallègue et al., 2022). Additionally, the current study concentrated on the wind

speed and direction at both 10-m and 50-m heights for precipitation forecasting due to their vital importance in evaluating evaporation, surface heat fluxes, momentum transfer, and atmospheric dynamics, which are particularly significant in areas with diverse terrains. The entire data can be categorised into six primary parameters for better understanding:

- Precipitation parameter: Rainfall.
- Temperature parameters: T2M, T2MDEW, T2MWET, TS, T2M_RANGE, T2M_MAX, and T2M_MIN.
- Humidity parameters: QV2M and RH2M.
- Pressure parameter: PS.
- Wind speed parameters: WS10M, WS10M_MAX, WS10M_MIN, WS10M_RANGE, WS50M, WS50M_MAX, WS50M_MIN, and WS50M_RANGE.
- Wind direction parameters: WD10M and WD50M.

Standard preprocessing steps were applied, beginning with the integration of the year, month, and day fields into a single “Date” column to streamline the temporal data format. Redundant original separate columns were then removed. Subsequently, a comprehensive examination of the dataset was conducted to identify and address issues such as duplicates, null entries, missing values, and outliers.

4.2. Exploratory data analysis

The analysis of rainfall patterns across the four regions reveals several salient trends. When examining the monthly distribution of rainfall, a conspicuous pattern emerges wherein England experiences its most substantial rainfall during the months of October, November, December, and January (Fig. 7). This noteworthy trend is consistently

Table 1

Description of data fields.

Variable names	Description
PS	Surface Pressure (kPa)
WS10M	Wind Speed at 10 Metres (m/s)
WS10M_MAX	Wind Speed at 10 Metres Maximum (m/s)
WS10M_MIN	Wind Speed at 10 Metres Minimum (m/s)
WS10M_RANGE	Wind Speed at 10 Metres Range (m/s)
WD10M	Wind Direction at 10 Metres (Degrees)
QV2M	Specific Humidity at 2 Metres (g/kg)
RH2M	Relative Humidity at 2 Metres (%)
PRECTOTCORR	Precipitation Corrected (mm/day)
T2M	Temperature at 2 Metres (C)
T2MDEW	Dew/Frost Point at 2 Metres (C)
T2MWET	Wet Bulb Temperature at 2 Metres (C)
TS	Earth Skin Temperature (C)
T2M_RANGE	Temperature at 2 Metres Range (C)
T2M_MAX	Temperature at 2 Metres Maximum (C)
T2M_MIN	Temperature at 2 Metres Minimum (C)
WD50M	Wind Direction at 50 Metres (Degrees)
WS50M_RANGE	Wind Speed at 50 Metres Range (m/s)
WS50M_MIN	Wind Speed at 50 Metres Minimum (m/s)
WS50M_MAX	Wind Speed at 50 Metres Maximum (m/s)
WS50M	Wind Speed at 50 Metres (m/s)

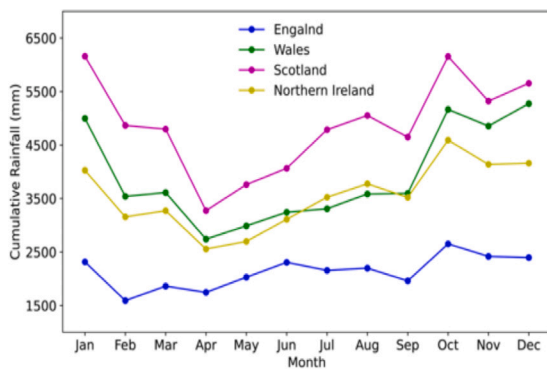


Fig. 7. Monthly average of rainfall.

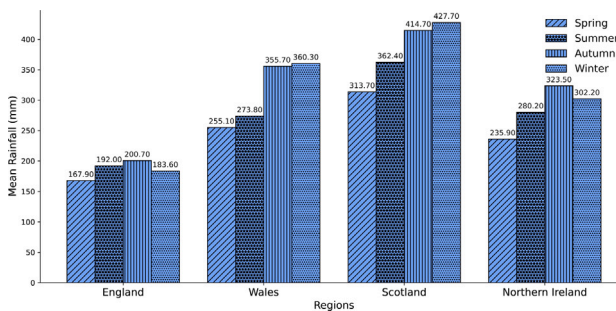


Fig. 8. Seasonal rainfall of 4 regions in 1981-2023.

observed in the other three regions as well. In terms of monthly averages, all regions, except for England, exhibit a consistent pattern, with October consistently being the wettest month and April consistently being the driest throughout the entire study period. England has the lowest monthly average during the February months, indicating a spatial heterogeneity influence on the dataset.

Fig. 8 depicts the seasonal rainfall patterns across England, Wales, Scotland, and Northern Ireland. England consistently experiences its highest rainfall during autumn, closely followed by summer. Conversely, Wales and Scotland exhibit peak rainfall during winter, with autumn closely trailing. Scotland consistently has the highest rainfall across all seasons among the four regions. Northern Ireland witnesses substantial rainfall during autumn, followed by winter. Overall, the

transition from mid-autumn to mid-winter emerges as the period of peak rainfall across all regions, highlighting a consistent pattern.

Our investigation into seasonality’s influence on rainfall patterns aligns with previous studies (Praveen et al., 2020; Manandhar et al., 2019; Zahran et al., 2023). These studies emphasise the significance of incorporating seasonal variations into rainfall prediction models. Recognising seasonality as a pivotal factor in the UK’s rainfall dynamics, we have integrated it as a feature in our dataset.

A coherent wave-shaped trend can be observed in the yearly cumulative rainfall across all regions as shown in Fig. 9. This trend is characterised by gradual increases over several years, reaching specific peaks, and then gradually declining in subsequent years. Notably, a distinct decreasing trend in cumulative rainfall has been observed in all regions since 2020. The gap between the latest annual values across the regions varies significantly from the annual average of the last ten years, with the Wales region showing the highest gap among them.

A detailed analysis was conducted to investigate the impact of temperature on rainfall and humidity by examining the percentage increase observed in the first and last five years of the period 1981-2023 for these regions as shown in Fig. 10. Scotland, the region with the lowest temperatures, exhibited the highest increase in temperature at 19.49%, surpassing the other regions, with England following closely. Surprisingly, this rise in temperature in the Scottish region had an inverse effect on rainfall, resulting in a lower increase in rainfall in that region. Conversely, regions with higher temperatures, such as England and Wales, experienced the most significant increases in rainfall percentages. The increase in temperature had varying effects on different regions. Specifically, Scotland, the coldest among the four regions, witnessed a substantial rise in temperature, but this rise had an inversely proportional impact on rainfall, leading to a lower increase. In contrast, the warmer regions, England and Wales, experienced the highest increases in rainfall. Furthermore, the rise in temperature had a direct influence on humidity in the Scotland region, leading to the highest percentage increase in humidity. In England, in contrast, the increase in temperature was associated with a rise in rainfall but the lowest increase in humidity. The UK as a whole region is experiencing a substantial rise in temperature, ranging from 10% to 20%. However, the impact of this temperature rise varies significantly depending on the geographical location.

4.3. Feature selection

Feature selection in rainfall forecasting serves as a validated technique to diminish dimensionality and augment model efficacy by pinpointing and retaining the most influential meteorological variables (Roster et al., 2022; Zhang et al., 2023). This method has proved to mitigate overfitting and enhance computational efficiency by eliminating redundant features, ensuring that deep learning models, particularly LSTM, operate optimally (Barrera-Animas et al., 2022; Sun et al., 2020). Guided by the correlation matrix, our selection process prioritises critical variables like temperature and humidity, pivotal for accurate rainfall prediction (Kim et al., 2023). The correlation matrix results for the England region dataset are depicted in Fig. 11. During feature selection, feature pairs with a correlation value of ± 0.7 or less are retained for each region individually (Barrera-Animas et al., 2022), while features exhibiting high correlations across all four regions are assessed for potential removal.

Our analysis reveals the following:

- Temperature-related features emphasise the primary temperature feature (T2M) as the most relevant. T2M is retained, while redundant features, such as those containing maximum, minimum, earth skin, wet temperature, dew point, and temperature range values, are excluded (Quej et al., 2022).
- Humidity-related measurements exhibit notable correlations. The primary humidity feature (RH2M) is retained, and the specific humidity feature is eliminated.

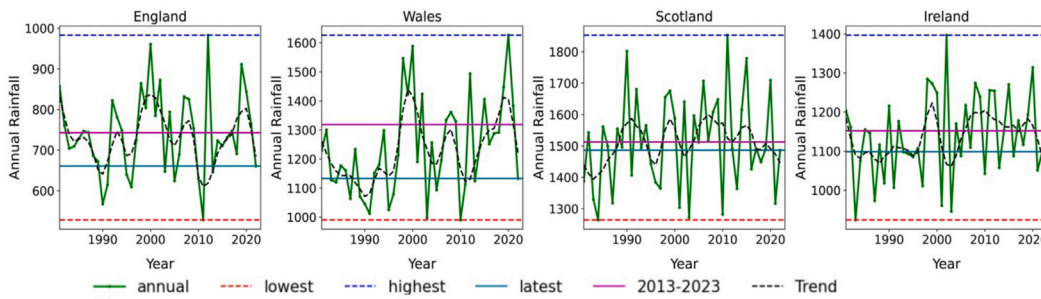


Fig. 9. Annual rainfall analysis.

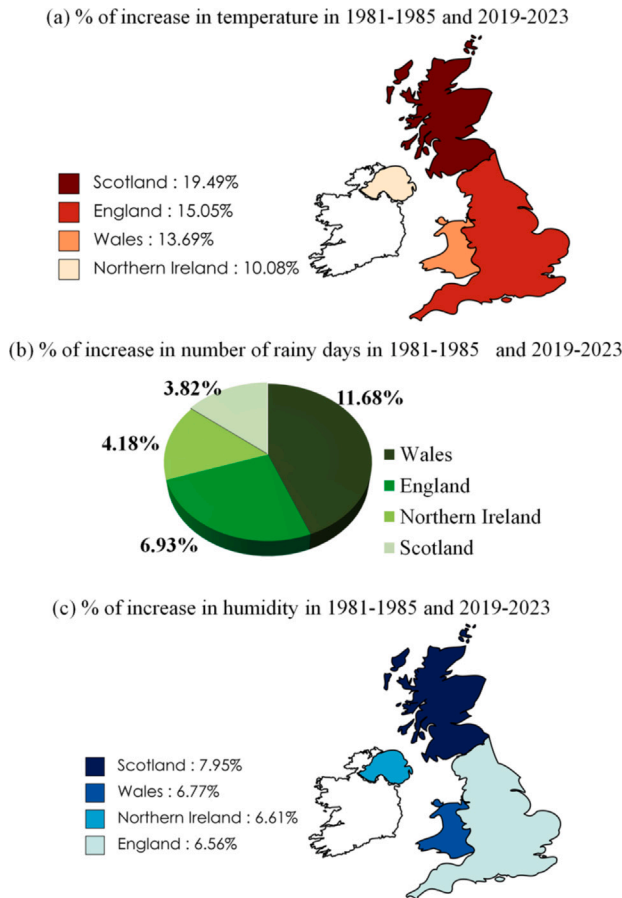


Fig. 10. Impact of temperature increase on rainfall and humidity from 1981–2023.

- Wind speed measurements show significant correlations. The main wind speed feature (WS10M) is retained, while related features are excluded.
- Wind direction features display strong correlations. The primary wind direction feature (WD10M) is retained, and WD50M is excluded.

The final dataset after feature engineering contains date, rainfall, surface pressure, wind speed (10M), wind direction (10M), relative humidity (2M), temperature (2M), and season. Among them all features except date fed into the deep learning model for forecasting.

To address variations in feature values, dataset normalisation is applied alongside feature selection (Kim and Bae, 2017). Specifically,

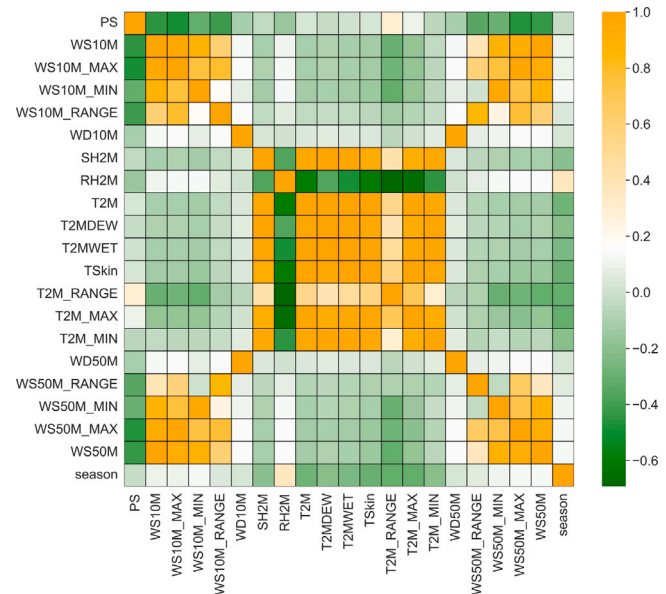


Fig. 11. Correlation matrix of England dataset.

the MinMaxScaler¹ has been used for data normalisation, transforming features to a 0–1 range. Normalisation is particularly critical when implementing the LSTM model to preserve data distribution and maintain consistent scales for numerical columns (Shanker et al., 1996). These structured datasets, following pre-processing, feature engineering, and normalisation, are employed for training and testing the rainfall prediction models.

4.4. Modelling

LSTM-based models

The first goal of this study is to adapt classic LSTM-based models for predicting daily rainfall. We consider three LSTM-based models as benchmarks: LSTM, Stacked-LSTM, and Bidirectional LSTM Networks. Since the effectiveness of these machine learning algorithms heavily depends on parameter and hyperparameter selection, we conducted a hyperparameter grid search following recommendations from Young et al. (2015).

To systematically explore the hyperparameter space and identify optimal configurations for each model, we employed a grid search technique. Key hyperparameters such as time steps, activation functions, number of hidden units, learning rate, optimiser functions, number of hidden layers, and number of epochs were meticulously varied across

¹ The formula used by MinMaxScaler to scale each feature x is: $x_{scaled} = \frac{x - \min(x)}{\max(x) - \min(x)}$.

Table 2
Configuration of each tuned model.

LSTM model				
Parameter	England	Wales	Scotland	Northern Ireland
Time steps	4	4	4	4
Activation function	swish	ReLU	ReLU	swish
Units	128	128	128	128
Learning rate	0.001	0.01	0.1	0.1
Optimiser	SGD	SGD	SGD	SGD
Number of hidden layers	1	1	1	1
Number of epochs	50	50	50	50
Stacked LSTM model				
Parameter	England	Wales	Scotland	Northern Ireland
Time steps	4	4	4	4
Activation function	tanh	ReLU	tanh	ReLU
Units	64	64	64	64
Learning rate	0.001	0.001	0.001	0.1
Optimiser	adam	SGD	SGD	SGD
Number of hidden layers	1	1	1	1
Number of epochs	50	50	50	50
Bidirectional LSTM model				
Parameter	England	Wales	Scotland	Northern Ireland
Time steps	4	4	4	4
Activation function	swish	tanh	tanh	tanh
Units	64	64	128	128
Learning rate	0.001	0.01	0.001	0.001
Optimiser	adam	SGD	SGD	SGD
Number of hidden layers	1	1	1	1
Number of epochs	50	50	50	50

predefined ranges (Kumar et al., 2019; Kim and Bae, 2017; Aswin et al., 2018). For instance, learning rates were tested over logarithmically spaced values ranging from 10^{-4} to 10^{-1} , while batch sizes were explored from 32 to 128 samples per batch. The number of hidden units in the LSTM layers was varied from 64 to 256 units. Furthermore, various activation functions including ReLU, tanh, and sigmoid were evaluated to assess their impact on model performance. Each combination of hyperparameters was exhaustively evaluated by dividing the data into training, validation, and test sets, accounting for 67%, 17%, and 16% of temporal data, respectively.

The grid search process was computationally intensive but essential for identifying configurations that minimise the loss function across different datasets and model architectures. The significance of these hyperparameters was evident through observed performance variations during testing. Minor adjustments, such as modifying batch size or optimiser, notably affected LSTM network performance. Following this tuning phase, we successfully identified the optimal hyperparameters for each dataset, as illustrated in Table 2.

Hybrid models

In addition to the three LSTM-based models serving as a baseline, we introduce two hybrid models aimed at potentially enhancing baseline performance.

The initial proposed hybrid model, illustrated in Fig. 12, combines CNN with LSTM layers – abbreviated as CLSTM – to effectively process sequential data, particularly in time series forecasting. CNNs excel at extracting spatial features from input time series, which are then used by LSTM to capture temporal dependencies. This approach is well-suited for tasks where both local patterns and long-term dependencies are crucial.

The model commences with a 1D convolutional layer featuring 10 filters and a tanh activation function with a kernel size of 3, tailored for localised pattern recognition. Following this, a max-pooling layer with a pool size of 2 is employed to reduce spatial dimensions while preserving crucial information.

Two LSTM layers are incorporated into the model. The first LSTM layer comprises 128 units with a sigmoid activation function configured

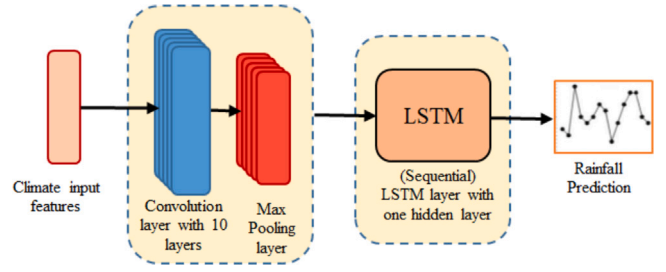


Fig. 12. Proposed CLSTM.

to return sequences, enabling it to capture long-term dependencies. To mitigate overfitting, a dropout layer with a rate of 0.2 follows this LSTM layer. A second hidden LSTM layer follows, comprising 64 units with a sigmoid activation function, incorporating a dropout layer with a 0.2 rate for regularisation. The model culminates with a dense layer serving as the output layer for regression tasks, consisting of a single neuron with a tanh activation function. During training, the Mean Squared Error (MSE) loss function is employed, with optimisation conducted using the Adam optimiser set to a learning rate of either 0.01 or 0.001, contingent on the dataset's characteristics within their respective regions. The selection of optimiser, activation functions, learning rate, and epoch settings is informed by insights gleaned from a review of existing models.

The data is first fed into the CNN network, producing an output as described in Eq. (1). The CLSTM method can be represented using Eqs. (1), (2), and (4) as follows Kim and Cho (2019):

$$\begin{aligned}
 \text{input gate } (i) &= (W_i P_t + h_{t-1} + U_{iX_t} + V_{iC_{t-1}}) \\
 \text{output gate } (o) &= (W_o P_t + h_{t-1} + U_{oX_t} + V_{oC_{t-1}}) \\
 \text{forget gate } (f) &= (W_f P_t + h_{t-1} + U_{fX_t} + V_{fC_{t-1}})
 \end{aligned} \tag{6}$$

$$\text{internal hidden state } (g) = \tanh(W_g h_{t-1} + U_g X_t)$$

$$\text{current cell state } c_t = (f * c_{t-1}) + (g * i)$$

$$\text{current hidden state } h_t = \tanh(c_t) * o$$

where P_t is the value of the pooling layer. The predicted rainfall y_t is given by:

$$y_t = f(W_0 * h_t)$$

The second hybrid model proposed, depicted in Fig. 13 and abbreviated as RLSTM, combines RNN and LSTM layers to leverage the unique advantages of both architectures. The RLSTM model capitalises on the sequential processing capabilities of RNNs, making it effective in capturing sequential patterns. Meanwhile, LSTMs address issues like vanishing gradients and excel at capturing long-term dependencies. This model is particularly well-suited for tasks involving input data with sequential dependencies, requiring the capture of both short-term and long-term patterns.

The architecture initiates with a Simple RNN layer comprising 128 units utilising the sigmoid activation function. Following this, a dropout layer with a 20% dropout rate is introduced for regularisation. Subsequently, an LSTM layer is integrated into the RNN layer, featuring 128 units with a sigmoid activation function aimed at capturing long-range dependencies within the sequential data. Another dropout layer, maintaining a 20% dropout rate, follows for further regularisation. Finally, a dense layer with a tanh activation function and a single unit is added to produce the model's output. During training, the Adam optimiser is employed with a learning rate set to 0.001, with the objective of minimising the MSE loss.

As an initial step, the data is fed into RNN network, and obtain the output O_t^R by using Eq. (7).

$$O_t^R = f(W * h_t + b) \tag{7}$$

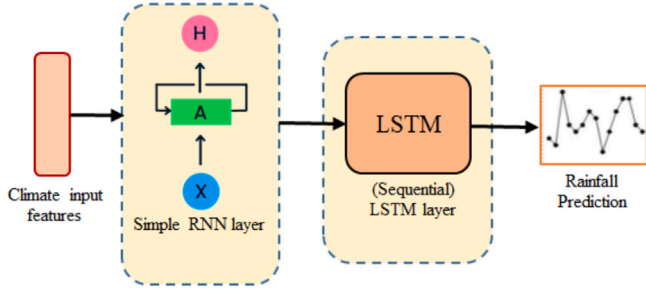


Fig. 13. Proposed RLSTM.

Table 3
Configuration of the two hybrid models.

CLSTM model				
Parameter	England	Wales	Scotland	Northern Ireland
Time steps	30	30	30	30
Activation function	tanh	tanh	tanh	sigmoid
Units	64	64	64	64
Learning rate	0.001	0.001	0.001	0.1
Optimiser	adam	adam	adam	adam
Number of hidden layers	1	1	1	1
Number of epochs	100	100	100	100
RLSTM model				
Parameter	England	Wales	Scotland	Northern Ireland
Time steps	15	30	30	15
Activation function	tanh	tanh	tanh	sigmoid
Units	64	64	64	64
Learning rate	0.001	0.001	0.001	0.1
Optimiser	adam	adam	adam	adam
Number of hidden layers	1	1	1	1
Number of epochs	100	100	100	100

O_t^R is the output/predicted value from RNN, b is the bias term and W is the weight; h_t is the hidden state for the current input value X_t at current time t . The O_t^R is then pass to the LSTM layer to get the final prediction y .

$$\begin{aligned}
 \text{input gate } (i) &= (W_i O_t^R + U_i X_t + V_i C_{t-1}) \\
 \text{output gate } (o) &= (W_o O_t^R + U_o X_t + V_o C_{t-1}) \\
 \text{forget gate } (f) &= (W_f O_t^R + U_f X_t + V_f C_{t-1}) \\
 \text{internal hidden state } (g) &= \tanh(W_g O_t^R + U_g X_t) \\
 \text{current cell state } c_t &= (f * c_{t-1}) + (g * i) \\
 \text{current hidden state } h_t &= \tanh(c_t) * o
 \end{aligned} \tag{8}$$

And finally, the predicted rainfall y_t is given by:

$$y_t = f(W_0 * h_t) \tag{9}$$

where h_t is the current hidden state and W_0 is the weight matrix of the output gate.

It is worth noting that the choice of optimiser, activation functions, learning rate, and the number of training epochs were adapted from the findings in the literature, as discussed in the previous section. This proposed hybrid architecture effectively leverages the capabilities of both RNN and LSTM layers to model both short-term and long-term dependencies in sequential data while mitigating overfitting through dropout regularisation.

Table 3 presents the optimal hyperparameter configurations for the two hybrid models across each region of the UK.

Table 4
Performance of each tuned model.

LSTM-based models				
Model	Region	Loss	RMSE	MAE
LSTM	England	0.013	0.114	0.068
	Wales	0.025	0.157	0.095
	Scotland	0.015	0.121	0.083
	Ireland	0.012	0.108	0.071
Stacked LSTM	England	0.011	0.107	0.063
	Wales	0.028	0.167	0.105
	Scotland	0.015	0.121	0.086
	Ireland	0.012	0.108	0.075
Bidirectional LSTM	England	0.017	0.130	0.074
	Wales	0.025	0.158	0.099
	Scotland	0.013	0.116	0.077
	Ireland	0.012	0.111	0.077
Hybrid models				
Model	Region	Loss	RMSE	MAE
CLSTM	England	0.012	0.107	0.063
	Wales	0.024	0.154	0.093
	Scotland	0.014	0.119	0.079
	Ireland	0.011	0.107	0.070
RLSTM	England	0.011	0.107	0.062
	Wales	0.022	0.150	0.088
	Scotland	0.013	0.115	0.077
	Ireland	0.011	0.104	0.068

4.5. Accuracy matrix

The following metrics were used to evaluate the effectiveness of the trained rainfall prediction models:

- **Loss**: this metric measures the error rate of the model in producing accurate results. It is calculated as follows:

$$Loss = \begin{cases} 1 & \text{if error occurs} \\ 0 & \text{otherwise} \end{cases}$$

- **RMSE**: it evaluates the square root of the average of the squared differences between the model's predictions and the actual values. The formula for RMSE is:

$$RMSE_{fo} = \sqrt{\frac{\sum_{i=1}^n (z_{fi} - z_{oi})^2}{n}}$$

where f represents model outputs, o represents observations, and n is the sample size.

- **MAE**: it assesses the average of the absolute differences between the model's predictions and the actual values. Its formula is:

$$MAE_{fo} = \frac{\sum_{i=1}^n |z_{fi} - z_{oi}|}{n}$$

By using these metrics together, we can assess and compare the accuracy and performance of the trained rainfall prediction models.

5. Results

We segment this section into distinct parts. Initially, we computed the accuracy measures of all tuned models, including the three LSTM-based models and the two proposed hybrid models, using the unseen test data to quantitatively evaluate their ability to fit unseen rainfall data. Following this, we visualised the training and validation loss for different epochs to qualitatively assess whether our models are exhibiting signs of overfitting or underfitting. Finally, we conducted another visualisation to depict the actual and predicted rainfall values for each model, aiming to identify patterns or trends in their performance. This analysis helps discern whether certain models excel under specific weather conditions or in particular regions, while others may struggle to accurately predict rainfall across various scenarios.

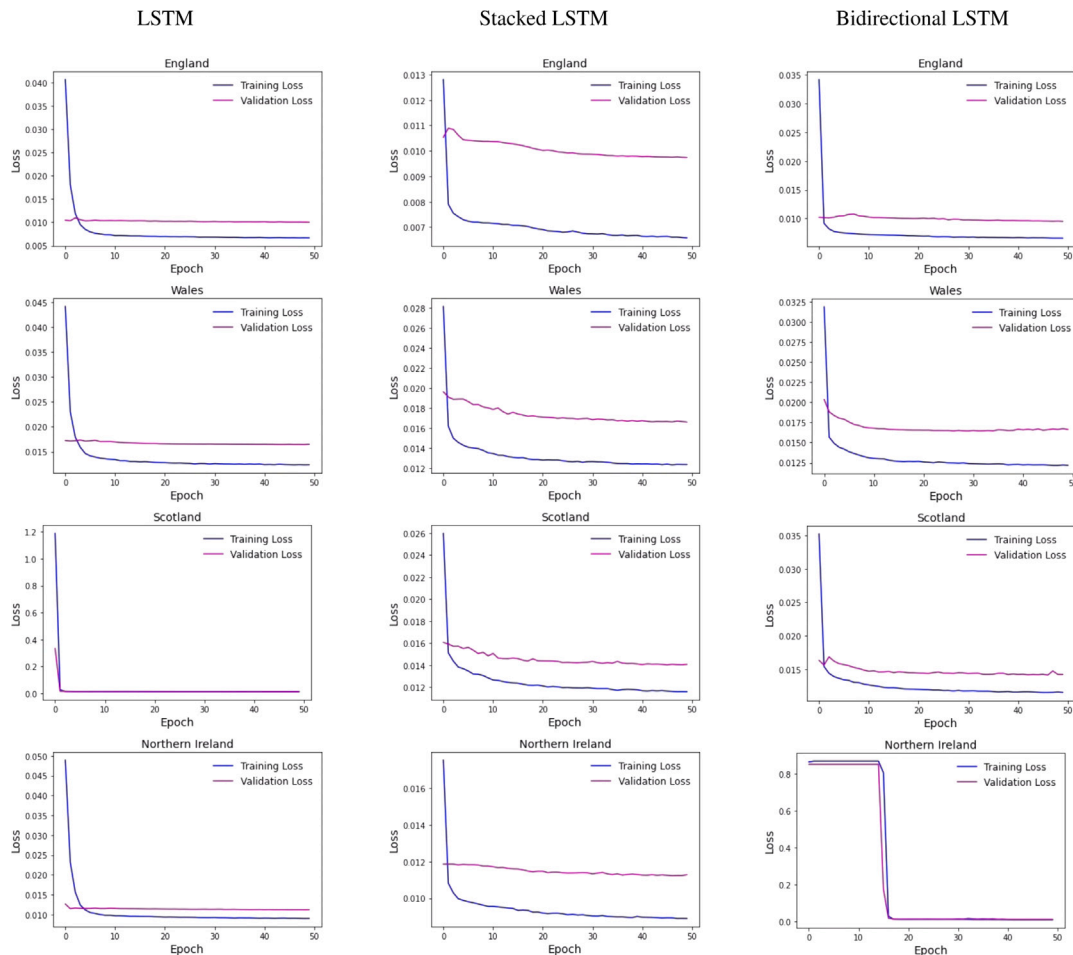


Fig. 14. Train and validation loss of LSTM-based models.

5.1. Accuracy measures of all tuned models

Table 4 provides an overview of the training outcomes for all five tuned models using unseen rainfall data. Across all four regions, the proposed hybrid RLSTM consistently outperforms the other models, boasting lower loss, RMSE, and MAE metrics. However, results reveal region-specific performance disparities: specifically, in the England region, the stacked LSTM model performs comparably to the RLSTM, surpassing other models. Similarly, the Bidirectional LSTM model matches the RLSTM's performance in the Northern Ireland region. For Wales and Scotland, the three LSTM-based models exhibit comparable loss values. Notably, both Bidirectional LSTM and Stacked LSTM excel in the Wales and Northern Ireland regions in terms of MAE values, while Bidirectional LSTM achieves the lowest MAE values for England and Scotland.

5.2. Training and validation loss

Figs. 14 and 15 illustrate the training and validation loss plots for LSTM-based and hybrid models, respectively. These plots serve as visual indicators of the model's performance throughout the training process, allowing for the detection of potential issues like overfitting, underfitting, and dataset representativeness. The objective of this analysis is to identify an optimal learning curve where the training and validation loss curves converge to a stable point while minimising the generalisation gap (Goodfellow et al., 2016).

Fig. 14 shows that classic LSTM model consistently achieve the best performance across all four regions of the UK, with the lowest training and validation losses and a stable convergence with a minimal gap between the curves. In contrast, the Stacked LSTM and Bidirectional LSTM models show larger gaps between their training and validation losses, indicating potential overfitting, particularly noticeable in the case of Stacked LSTM. Bidirectional LSTM performs relatively well with the Northern Ireland dataset, initially displaying a small training and validation gap that stabilises after 15 epochs. The observed differences in performance among the models suggest that LSTM is a more suitable choice for rainfall prediction, given its superior handling of overfitting (see Fig. 14).

Fig. 15 illustrates that RLSTM generally surpasses CLSTM in minimising the generalisation gap, except in the England region where the gap is marginally larger for the RLSTM model. RLSTM demonstrates remarkable performance, with training and validation loss curves converging to a stable point, especially in the Wales region.

5.3. Actual and predicted rainfall values

Figs. 16 and 17 illustrate the comparison between actual and predicted rainfall values generated by each model, focusing on a selected portion of unseen data. These visualisations aid in identifying models that perform well under specific weather conditions or in particular regions, while also highlighting models that may face challenges in accurately predicting rainfall across different scenarios.

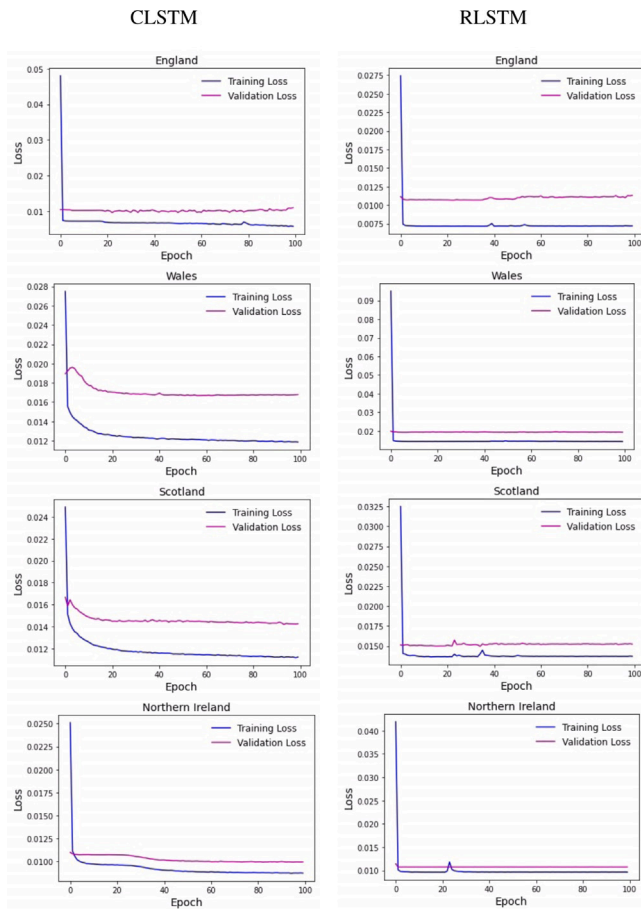


Fig. 15. Train and validation loss of hybrid models.

Fig. 16 shows that the classic LSTM and Bidirectional LSTM models manage to provide nice fitting predicted values for the Scotland and Wales regions and are also able to capture some peak values. When considering Stacked LSTM, it is able to capture some peaks in the data with minimal gaps between the predicted and actual rainfall data, especially for the Scotland region.

Fig. 17 focuses on the hybrid models instead. It shows that the CLSTM model predicts rainfall with minimal errors and performs better for rainfall prediction in the Wales and Scotland regions. However, among the five models, the RLSTM predicts the values more accurately, especially for the Wales and Scotland regions. The RLSTM captures some of the peak values of all four regions, especially in Scotland, England, and Wales regions.

The comparative analysis of Figs. 16 and 17 reveals that, overall, the CLSTM and RLSTM models outperform the three LSTM-based models in predicting rainfall across all four regions of the UK using meteorological data. However, among all tested models, RLSTM stands out as the top performer, surpassing even the CLSTM model. Specifically:

- RLSTM demonstrates great capabilities in capturing sequential dependencies inherent in meteorological time series data, a critical aspect in rainfall forecasting tasks that rely on historical weather patterns to predict future rainfall. We speculate that its effectiveness lies in its ability to model the temporal evolution of meteorological variables over time intervals, making it particularly suitable for the dynamic and intricate weather patterns observed in all UK regions, highlighting the importance of considering temporal aspects when modelling weather patterns.

- In contrast, the CLSTM model emphasises the extraction of spatial features from meteorological data using CNNs before employing LSTM networks for temporal modelling. While CNNs are adept at capturing spatial patterns such as temperature distributions or wind patterns, they may not comprehensively capture the complex temporal dependencies inherent in meteorological data, especially in regions characterised by highly variable weather conditions like those in the UK.

6. Discussion and conclusion

The research conducted in this study aims to fill existing gaps in rainfall forecasting methodologies, with the overarching objective of improving accuracy and deepening our understanding of temperature variations' impact on rainfall patterns across all regions of the UK, including England, Wales, Northern Ireland, and Scotland.

This study offers several contribution, specifically:

- It optimise three distinct LSTM-based models – LSTM, Stacked-LSTM, and Bidirectional LSTM Networks – using time-series data from the four UK countries to predict daily rainfall.
- It introduces two hybrid models – CNN with LSTM (CLSTM) and RNN with LSTM (RLSTM) – for predicting daily rainfall using time-series data from the four UK countries.
- It evaluates the performance of each model in terms of their ability to forecast daily rainfall amounts using time-series data from the four UK countries.

The proposed rainfall forecasting models present invaluable insights and predictions that can profoundly influence decision-making processes in agriculture, water resource management, and disaster preparedness. Through accurate and timely forecasts, these models empower farmers to make informed decisions regarding planting schedules, irrigation management, and crop protection measures, thereby optimising agricultural productivity. Furthermore, water resource managers can utilise these forecasts to more effectively allocate water resources, mitigate the risks of floods or droughts, and plan infrastructure projects. Additionally, timely predictions support disaster preparedness efforts by enabling authorities to implement proactive measures such as evacuation plans, emergency response strategies, and resource allocation, thereby minimising the impact of extreme weather events on communities and infrastructure. Overall, the adoption of these forecasting models enhances resilience and fosters sustainable development across various sectors.

The results of the comprehensive UK rainfall forecasting study reveal a complex interplay between meteorological parameters, deep learning models, and regional rainfall patterns. Notably, temperature affects rainfall differently across UK regions, with Scotland experiencing decreasing rainfall despite rising temperatures, while England and Wales see increased precipitation. These findings underscore the importance of considering regional disparities in model selection.

Among the five machine learning models evaluated, the hybrid RLSTM consistently outperformed others in all regions across various evaluation metrics, including loss, MAE, and RMSE. Visual analysis of training and validation loss curves revealed that classic LSTM performed well in all regions, while Stacked LSTM produced less accurate results. Bidirectional LSTM exhibited potential overfitting in most regions. Both RLSTM and CLSTM, incorporating hybrid LSTM networks, displayed favourable loss curves, with RLSTM excelling in minimising the generalisation gap. However, the comparison of actual and predicted rainfall values highlighted the challenge of forecasting peak rainfall events. Although CLSTM and RLSTM showed promise, particularly in Wales and Scotland, they encountered difficulties in forecasting sudden precipitation variations, suggesting the need for further refinement.

We posit that the challenge of forecasting peak rainfall events stems from the absence of crucial features in the existing data. Parameters

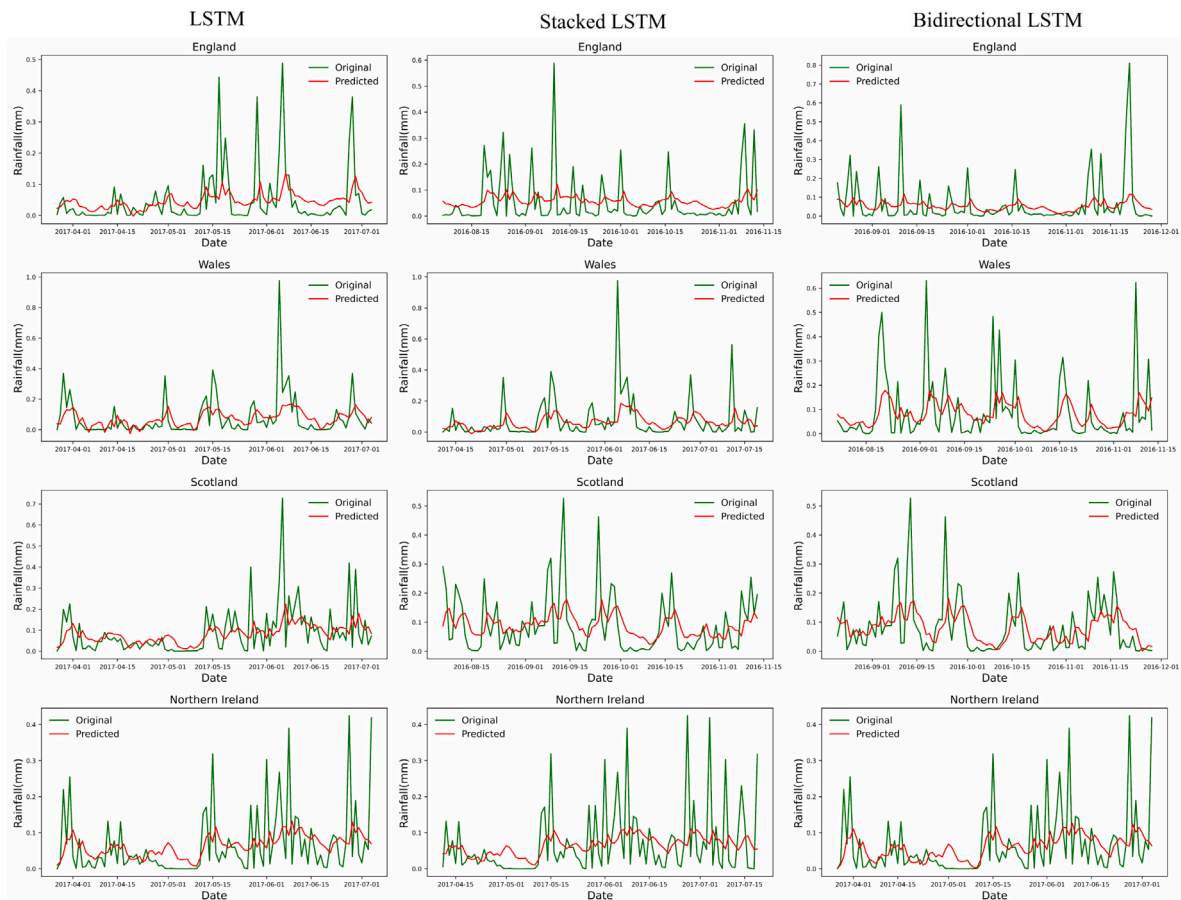


Fig. 16. Actual and predicted rainfall on a portion of unseen data using LSTM models.

such as specific cloud characteristics, humidity levels, and snow presence could play pivotal roles in precise rainfall prediction (Hemavathi et al., 2021; SS, 2023). These components are integral to the intricate climate system, particularly in regions like the UK, where weather patterns are highly variable and influenced by diverse meteorological factors. Incorporating multi-modal data, including vegetation cover and other environmental variables that significantly impact rainfall patterns (Wu and Li, 2023), could be essential for improving predictions. Hence, addressing these data gaps and enhancing data representation could lead to more precise rainfall forecasting in the future.

While these findings hold promise, they also highlight the challenge of generalisation in LSTM-based models. To address this challenge, future research could explore the integration of uncertainty measures into predictions. Techniques such as bootstrap methods or deep quantile regression could serve as viable tools for quantifying the inherent uncertainty and variability in rainfall prediction, thereby enhancing the reliability and interpretability of model predictions. Additionally, leveraging hourly data and incorporating additional meteorological variables, such as snow and cloud data, could further improve model performance. Efforts should also be directed towards developing models capable of accurately predicting peak rainfall values, which are often underestimated by current deep learning models. Ultimately, future research should prioritise enhancing forecast accuracy and exploring the application of models in flood and drought warning systems, thereby bolstering resilience to climate-related hazards, safeguarding vital resources, and fostering sustainable development efforts worldwide.

CRediT authorship contribution statement

Geethu Thottungal Harilal: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Resources, Project

administration, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Aniket Dixit:** Writing – review & editing, Supervision. **Giovanni Quattrone:** Supervision, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

The raw data collected from the NASA Power Data Viewer is available for researchers and the materials that the authors used are available at the authors' hands.

Acknowledgements

We gratefully acknowledge the NASA Power Data Viewer system from where the entire dataset was collected for the completion of this study.

Funding

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

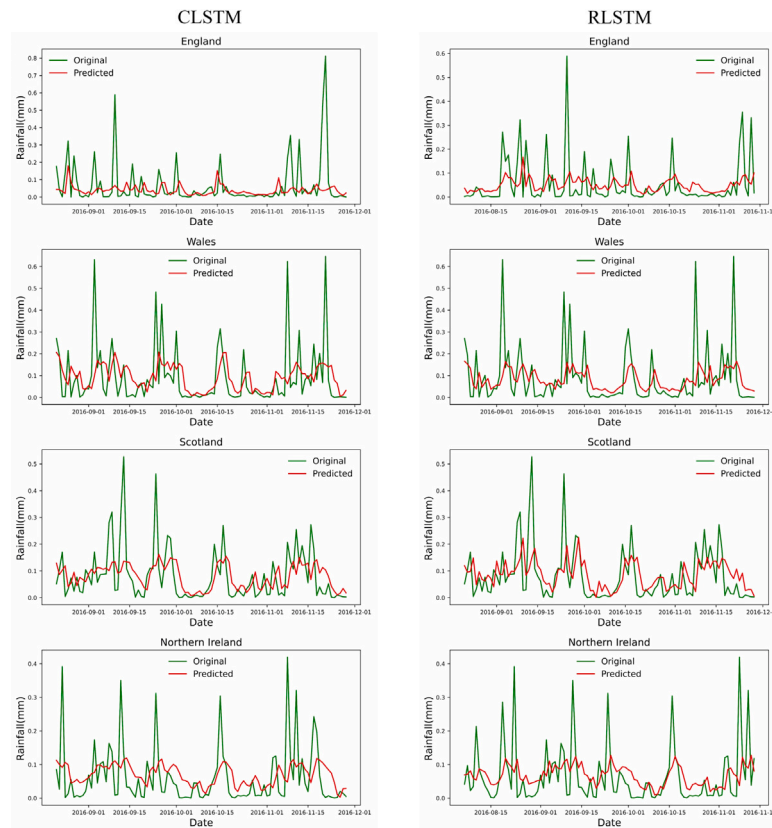


Fig. 17. Actual and predicted rainfall on a portion of unseen data using hybrid models.

References

- Aswin, S., Srikanth, D., Ravi, V., 2018. Deep learning models for the prediction of rainfall. pp. 0657–0661. <http://dx.doi.org/10.1109/ICCSP.2018.8523829>.
- Balluff, S., Bendfeld, J., Krauter, S., 2020. Meteorological data forecast using RNN. In: *Deep Learning and Neural Networks: Concepts, Methodologies, Tools, and Applications*. IGI Global, pp. 905–920.
- Barrera-Animas, A., Oyedele, L., Bilal, M., Akinosho, T., Delgado, J., Akanbi, L., 2022. Rainfall prediction: A comparative analysis of modern machine learning algorithms for time-series forecasting. *Mach. Learn. Appl.* 7, 100204.
- Ben Bouallègue, Z., Cooper, F., Chantry, M., Düben, P., Bechtold, P., Sandu, I., 2022. Statistical modelling of 2 m temperature and 10 m wind speed forecast errors. *Mon. Weather. Rev.*
- Chao, Z., Pu, F., Yin, Y., Han, B., Chen, X., 2018. Research on real-time local rainfall prediction based on MEMS sensors. *J. Sens.*
- Cheng, H., Xie, Z., Wu, L., Yu, Z., Li, R., 2019. Data prediction model in wireless sensor networks based on bidirectional LSTM. *EURASIP J. Wireless Commun. Networking* 1–12.
- Cotterill, D., Stott, P., Christidis, N., Kendon, E., 2021. Increase in the frequency of extreme daily precipitation in the United Kingdom in autumn. *Weather Clim. Extrem.* 33, 100340.
- Cui, Z., Ke, R., Pu, Z., Wang, Y., 2018. Deep bidirectional and unidirectional LSTM recurrent neural network for network-wide traffic speed prediction. *arXiv preprint arXiv:1801.02143*.
- Gan, K., Sun, S., Wang, S., Wei, Y., 2018. A secondary-decomposition-ensemble learning paradigm for forecasting PM_{2.5} concentration. *Atmos. Pollut. Res.* 9 (6), 989–999.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT Press.
- GOV.UK, 2023. Climate change explained. URL [https://www.gov.uk/guidance/climate-change-explained#:~:text=The%20United%20Kingdom%20\(UK\)%20is](https://www.gov.uk/guidance/climate-change-explained#:~:text=The%20United%20Kingdom%20(UK)%20is).
- Greff, K., Srivastava, R.K., Koutník, J., Steunebrink, B.R., Schmidhuber, J., 2016. LSTM: A search space odyssey. *IEEE Trans. Neural Netw. Learn. Syst.* 28 (10), 2222–2232.
- Haidar, A., Verma, B., 2016. A genetic algorithm based feature selection approach for rainfall forecasting in sugarcane areas. In: *2016 IEEE Symposium Series on Computational Intelligence*. SSCI, IEEE, pp. 1–8.
- He, R., Zhang, L., Chew, A.W.Z., 2022. Modeling and predicting rainfall time series using seasonal-trend decomposition and machine learning. *Knowl.-Based Syst.* 251, 109125.
- Hemavathi, D., Kumar, V.R., Regin, R., Rajest, S.S., Phasinam, K., Singh, S., 2021. Technical support for detection and prediction of rainfall. In: *2021 2nd International Conference on Smart Electronics and Communication*. ICOSEC, IEEE, pp. 1629–1634.
- Hernández, E., Sanchez-Anguix, V., Julian, V., Palanca, J., Duque, N., 2016. Rainfall prediction: A deep learning approach. In: *Hybrid Artificial Intelligent Systems: 11th International Conference, HAIS 2016*. Springer International Publishing, pp. 151–162.
- Hossain, I., Rasel, H.M., Imteaz, M.A., Mekanik, F., 2020. Long-term seasonal rainfall forecasting using linear and non-linear modelling approaches: a case study for western Australia. *Meteorol. Atmos. Phys.* 132, 131–141.
- Hussein, E.A., Ghaziagar, M., Thron, C., Vaccari, M., Jafta, Y., 2022. Rainfall prediction using machine learning models: Literature survey. In: *Artificial Intelligence for Data Science in Theory and Practice*. pp. 75–108.
- Hutter, F., Kotthoff, L., Vanschoren, J., 2019. *Automated Machine Learning: Methods, Systems, Challenges*. Springer Nature, p. 219.
- Kim, H.U., Bae, T.S., 2017. Preliminary study of deep learning-based precipitation prediction, 35 (5), 423–429.
- Kim, T.Y., Cho, S.B., 2019. Predicting residential energy consumption using CLSTM neural networks. *Energy* 182, 72–81.
- Kim, J.H., Lee, H., Byeon, S., Shin, J.K., Lee, D.H., Jang, J., Chon, K., Park, Y., 2023. Machine learning-based early warning level prediction for cyanobacterial blooms using environmental variable selection and data resampling. *Toxics* 11 (12), 955.
- Kratzert, F., Klotz, D., Brenner, C., Schulz, K., Herrnegger, M., 2018. Rainfall-runoff modelling using long short-term memory (LSTM) networks. *Hydrol. Earth Syst. Sci.* 22 (11), 6005–6022.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems*. Vol. 25.
- Kumar, D., Singh, A., Samui, P., Jha, R.K., 2019. Forecasting monthly precipitation using sequential modelling. *Hydrol. Sci. J.* 64 (6), 690–700.
- Liu, Q., Zou, Y., Liu, X., Linge, N., 2019. A survey on rainfall forecasting using artificial neural network. *Int. J. Embed. Syst.* 11 (2), 240–249.
- Manandhar, S., Dev, S., Lee, Y.H., Meng, Y.S., Winkler, S., 2019. A data-driven approach for accurate rainfall prediction. *IEEE Trans. Geosci. Remote Sens.* 57 (11), 9323–9331.
- Markuna, S., Kumar, P., Ali, R., Vishwakarma, D.K., Kushwaha, K.S., Kumar, R., Singh, V.K., Chaudhary, S., Kuriqi, A., 2023. Application of innovative machine learning techniques for long-term rainfall prediction. *Pure Appl. Geophys.* 180 (1), 335–363.
- Nakicenovic, N., Alcamo, J., Davis, G., Vries, B.D., Fenhann, J., Gaffin, S., Gregory, K., Grubler, A., Jung, T.Y., Kram, T., La Rovere, E.L., 2000. *Special report on emissions scenarios*.
- Neagh, L., Lomond, L., 2014. *United Kingdom of great britain and Northern Ireland*.

- Ni, L., Wang, D., Singh, V.P., Wu, J., Wang, Y., Tao, Y., Zhang, J., 2020. Streamflow and rainfall forecasting by two long short-term memory-based models. *J. Hydrol.* 583, 124296.
- Ojo, S.O., Owolawi, P.A., Mphahlele, M., Adisa, J.A., 2019. Stock market behaviour prediction using stacked LSTM networks. In: 2019 International Multidisciplinary Information Technology and Engineering Conference. IMITEC, IEEE, pp. 1–5.
- Parmar, A., Mistree, K., Sompura, M., 2017. Machine learning techniques for rainfall prediction: A review. In: International Conference on Innovations in Information Embedded and Communication Systems. Vol. 3.
- Poornima, S., Pushpalatha, M., 2019. Prediction of rainfall using intensified LSTM based recurrent neural network with weighted linear units. *Atmosphere* 10 (11), 668.
- Power, N.A.S.A., 2023. Data Access Viewerz. URL <https://power.larc.nasa.gov/data-access-viewer>.
- Praveen, B., Talukdar, S., Shahfahad, Mahato, S., Mondal, J., Sharma, P., Islam, A.R.M.T., Rahman, A., 2020. Analyzing trend and forecasting of rainfall changes in India using non-parametrical and machine learning approaches. *Sci. Rep.* 10 (1), 10342.
- Quej, V.H., de la Cruz Castillo, C., Almorox, J., Rivera-Hernandez, B., 2022. Evaluation of artificial intelligence models for daily prediction of reference evapotranspiration using temperature, rainfall and relative humidity in a warm sub-humid environment. *Italian Journal of Agrometeorology* (1), 49–63.
- Ramsundram, N., Sathya, S., Karthikeyan, S., 2016. Comparison of decision tree based rainfall prediction model with data driven model considering climatic variables. *Irrig. Drain. Syst. Eng.* 5 (3), 1–5.
- Roster, K., Connaughton, C., Rodrigues, F.A., 2022. Machine-learning-based forecasting of dengue fever in Brazilian cities using epidemiologic and meteorological variables. *Am. J. Epidemiol.* 191 (10), 1803–1812.
- Saikhu, A., Arifin, A.Z., Faticah, C., 2017. Rainfall forecasting by using autoregressive integrated moving average, single input and multi input transfer function. In: 2017 11th International Conference on Information & Communication Technology and System. ICTS, IEEE, pp. 85–90.
- Shanker, M., Hu, M.Y., Hung, M.S., 1996. Effect of data standardization on neural network training. *Omega* 24 (4), 385–397.
- Singh, P., Borah, B., 2013. Indian summer monsoon rainfall prediction using artificial neural network. *Stoch. Environ. Res. Risk Assess.* 27, 1585–1599.
- Singh, U., Chauhan, S., Krishnamachari, A., Vig, L., 2015. Ensemble of deep long short term memory networks for labelling origin of replication sequences. In: 2015 IEEE International Conference on Data Science and Advanced Analytics. DSAA, IEEE, pp. 1–7.
- SS, D., 2023. A novel model for rainfall prediction using hybrid stochastic-based Bayesian optimization algorithm. *Environ. Sci. Pollut. Res. Int.*
- Sun, S., Fu, J., Zhu, F., Du, D., 2020. A hybrid structure of an extreme learning machine combined with feature selection, signal decomposition and parameter optimization for short-term wind speed forecasting. *Trans. Inst. Meas. Control* 42 (1), 3–21.
- Trenberth, K.E., 2005. The impact of climate change and variability on heavy precipitation, floods, and droughts. In: *Encyclopedia of Hydrological Sciences*. Vol. 17, pp. 1–11.
- Venkata Ramana, R., Krishna, B., Kumar, S.R., Pandey, N.G., 2013. Monthly rainfall prediction using wavelet neural network analysis. *Water Resources Management* 27, 3697–3711.
- Wang, Z., Yan, W., Oates, T., 2017. Time series classification from scratch with deep neural networks: A strong baseline. In: 2017 International Joint Conference on Neural Networks. IJCNN, IEEE, pp. 1578–1585.
- Wu, C.L., Chau, K.W., 2013. Prediction of rainfall time series using modular soft computing methods. *Eng. Appl. Artif. Intell.* 26 (3), 997–1007.
- Wu, S., Li, X., 2023. Rainfall-runoff prediction based on multi-modal data fusion. In: International Conference on Algorithms, High Performance Computing, and Artificial Intelligence (AHPACAI 2023). Vol. 12941, SPIE, pp. 322–326.
- Wu, D., Wu, L., Zhang, T., Zhang, W., Huang, J., Wang, X., 2022. Short-term rainfall prediction based on radar echo using an improved self-attention PredRNN deep learning model. *Atmosphere* 13 (12), 1963.
- Xu, L., Chen, N., Zhang, X., Chen, Z., 2020. A data-driven multi-model ensemble for deterministic and probabilistic precipitation forecasting at seasonal scale. *Clim. Dyn.* 54, 3355–3374.
- Young, S.R., Rose, D.C., Karnowski, T.P., Lim, S.H., Patton, R.M., 2015. Optimizing deep learning hyper-parameters through an evolutionary algorithm. In: Proceedings of the Workshop on Machine Learning in High-Performance Computing Environments. pp. 1–5.
- Yunpeng, L., Di, H., Junpeng, B., Yong, Q., 2017. Multi-step ahead time series forecasting for different data patterns based on LSTM recurrent neural network. In: 2017 14th Web Information Systems and Applications Conference. WISA, IEEE, pp. 305–310.
- Zahran, B., Ayyoub, B., Abu-Ain, W., Hadi, W., Al-Hawary, S., 2023. A fuzzy based model for rainfall prediction. *Int. J. Data Netw. Sci.* 7 (1), 97–106.
- Zhang, L., Xue, Z., Liu, H., Li, H., 2023. Enhanced generalized regression neural network with backward sequential feature selection for machine-learning-driven soil moisture estimation: A case study over the Qinghai-Tibet Plateau. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*