

# Novel machine learning approach for solar photovoltaic energy output forecast using extra-terrestrial solar irradiance

Cornelia A. Fjelkestam Frederiksen, Zuansi Cai \*

*School of Engineering and the Built Environment, Edinburgh Napier University, Edinburgh EH10 5DT, Scotland, UK*

## HIGHLIGHTS

- Extra-terrestrial Solar Irradiance has been validated for PV output forecasting.
- The machine learning approach successfully captures huge intra-daily PV output variations.
- Study paves the way to develop a simple and effective PV output forecasting approach.

## ARTICLE INFO

### Keywords:

Photovoltaic power prediction  
Extra-terrestrial irradiance  
Machine learning  
Non-linear Autoregressive Exogenous Neural Network

## ABSTRACT

The inherently intermittent nature of solar irradiance and other meteorological variables means that accurate forecasting of the photovoltaic power output is essential for planning and balancing photovoltaic power systems. This study proposes a novel approach to predicting one-week-ahead half-hourly photovoltaic power output in the United Kingdom using sloped extra-terrestrial irradiance and weather data (e.g., cloud-cover and temperature) as input parameters. A Non-linear Autoregressive Exogenous Neural Network is trained on a three-year historical dataset from two photovoltaic plants in the United Kingdom with capacities of 53 and 103 MWP. The forecasting model captures huge intra-daily variations of photovoltaic output, which is particularly useful to balance the supply and demand of the electricity system. The result of the study validates the concept of using sloped extra-terrestrial irradiance as an input parameter and suggests that meteorological conditions will dictate the accuracy of predictions. Findings also indicate that the use of sloped extra-terrestrial irradiance in conjunction with cloud-cover presented the optimal combination of input parameters as these provided the simplest and most cost-effective model without reducing accuracy. The approach can have universal value as it only requires coordinates and weather data. There is now a strong imperative to use the model in other locations where the weather is more stable.

## 1. Introduction and problematic

In efforts to meet the UK's net-zero greenhouse gas (GHG) target by 2050 [1], solar photovoltaic (PV) power has been undergoing rapid growth and increased demand. For example, during a 12-month period up to 31 March 2021, 975 megawatts (MW) PV capacity was installed in the UK. This increase brought the UK's total installed PV capacity to over 14 gigawatts (GW) [2], placing the UK in 6th place internationally for renewable energy generation [3].

Although PV energy is paving the way for clean energy, the power output of PV systems is intermittent by nature and largely dependent on weather and climate [4]. This dependency creates challenges for its

optimal use in the UK, where weather is notably stochastic. This unique variability in weather is related to Britain's position on the planet – an island situated between the Atlantic Ocean and continental Europe. As a result, the UK is under an area where five air masses meet, creating weather fronts. The UK's proximity to the polar front jet stream (a high-altitude ribbon of fast-moving air) also contributes to the unsettled weather, due to frequent changes in pressure [5].

The continued growth of PV energy, combined with volatility in the UK's weather, creates operable challenges for the UK National Grid, particularly in the summer months [6]. This is due to increased supply and demand variability caused by the low demand and high levels of renewable generation. For instance, the all-time peak from PV power in

\* Corresponding author.

E-mail address: [z.cai@napier.ac.uk](mailto:z.cai@napier.ac.uk) (Z. Cai).

**Table 1**

Information for the sources of data.

Data	Source	Description	Time Interval	Locations for Data
PV	Sheffield Solar [19]	A collaborative PV live service between the UK National Grid and the University of Sheffield. The service provides reliable time-series data for all solar PV systems connected to the UK transmission network.	30 min	Richborough & Grimsby
Meteorological Variables	Met Office [20]	MIDAS: UK hourly weather observation data is available in the CEDA archive. The archive holds historical measurements for weather stations across the UK, and the data spans from 1875 to present.	60 min	Manston & Donna Nook (nearby weather stations)
Modelled Sloped ERAD	Solar Radiation and Daylight Models [17]	The model is based on the theory of solar geometry and requires the following inputs: date and time, as well as the PV's coordinates, orientation and tilt angle. Relevant software is provided in the source.	30 min	Richborough & Grimsby substations
Modelled Irradiance	Copernicus Atmosphere Monitoring Service [21]	CAMS solar radiation time series provide modelled solar irradiance based on weather conditions through satellite images. Only historical data is available through this service.	15 min	Richborough & Grimsby substations

the UK was recorded at 9.86 GW on 20th of April 2020 at 12:30, whereas the peak generation two days prior (18th of April 2020 at 13:30) was less than half that amount at 4.73 GW [7].

To balance the supply and demand of the electricity system, the National Grid needs to take more actions to curtail renewable generation. For example, in 2020, wind curtailment cost UK electricity users £274 million [8]. This statistical data shows that, despite the UK struggling to meet carbon emission targets, electricity generated from renewables is still being wasted. As such, there is now a strong imperative to develop innovative approaches to high frequency forecasting of PV generation (e.g., by minute or hourly) for *7-days-ahead*.

Various models have been used for PV output forecasting, and most of the present literature focuses on short-term (up to *three-days-ahead*) forecasting [9]. The models are commonly divided into persistence methods, physical models and statistical approaches [10]. Persistence methods adopt the idea that the current day's climate is equal to the prevailing conditions of the previous day and is used for short-term and very-short-term forecasting [10]. Physical models commonly use numerical atmospheric data to forecast weather [11] and are mainly developed using Numerical Weather Prediction (NWP) [12]. Physical models perform best when meteorological conditions are stable [13] and

are often used for longer forecast horizons. Statistical models utilise mathematical equations to extract patterns from input data. Statistical techniques can be divided into two groups: time-series and machine learning (ML) based models [10]. Time-series based models are often used for short-term forecasting and use past values through assessing the pattern of past information [10].

ML models are based on computing or artificial intelligence (AI). The models utilise AI's ability to learn from historical data patterns and improve prediction with training runs. Artificial Neural Networks (ANN) is considered the most successful method for PV forecasting and is used widely thanks to its ability to model non-linear, complex, and dynamic processes [10]. For example, Andersson & Yakimenko [14] explored a Non-linear Auto Regressive Exogenous Neural Network (NARXNN) to forecast the PV output for a microgrid. Liu, Liu, Sun, Li & Wennersten [15] adopted a Backpropagation Neural Network (BPNN) for *24-hour-ahead* solar PV prediction. The superiority of ML to other models is particularly evident for longer time horizons. However, training is more complicated for this type of approach, and large amounts of data are needed. Furthermore, there exists no commonly accepted superior way to construct a perfect model, making the process iterative.

Throughout the literature, most PV power forecasts using ML approaches rely on solar irradiance data. The relationship between solar irradiance and PV generation means that reliable solar irradiance data is needed to forecast PV output accurately [16]. However, to accurately predict local solar irradiance, it requires an accurate prediction of the weather and long-term measurement of solar irradiance. This is particularly challenging for the UK given that there are only three weather stations recording hourly direct irradiance [17], as the uncertainty of solar irradiance prediction has proven to be a dominant source of PV forecasting error.

Limited local coverage of irradiation measurements often prompts the application of modelled irradiance, and parametric irradiation models (based on time, geographical location, and weather) are widely used. Su, Batzelis & Pal [9] used modelled irradiance from the Copernicus Atmosphere Monitoring Service (CAMS) in conjunction with modelled weather data to forecast PV output *6-days-ahead*. Similarly, Andersson & Yakimenko [14] predicted PV output for *24-hours-ahead* by

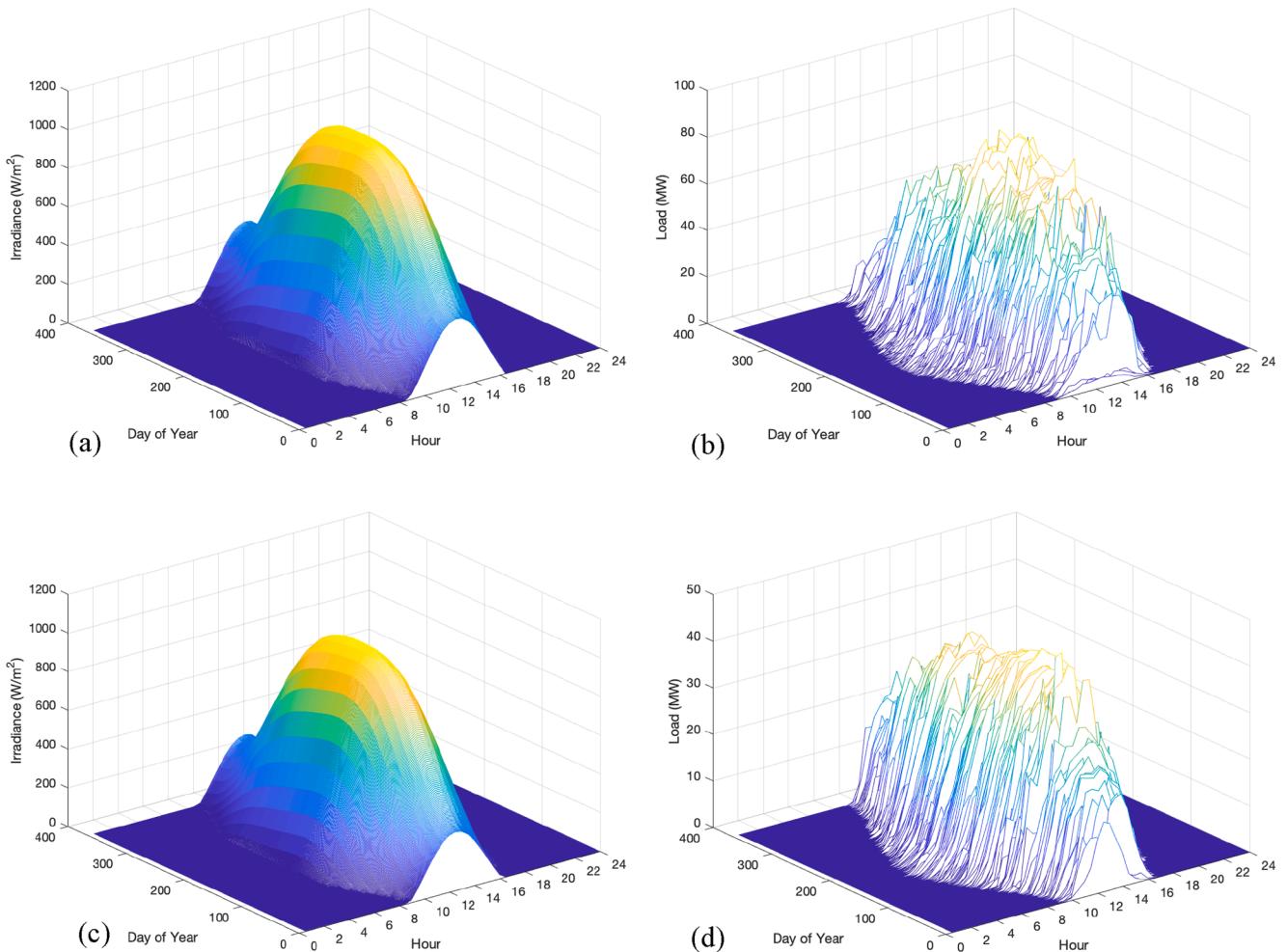


**Fig. 1.** Map illustrating locations in the UK. A denotes Grimsby substation & Donna Nook; and B denotes Richborough substation and Manston.

**Table 2**

Training and testing datasets.

Season	Training Data	Validation Data
Spring	01–03–2017 to 31–05–2017, 01–03–2018 to 31–05–2018, 01–03–2019 to 31–05–2019	19–03–2020 to 26–03–2020
Summer	01–06–2017 to 31–08–2017, 01–06–2018 to 31–08–2018, 01–06–2019 to 31–08–2019	28–07–2020 to 04–08–2020
Autumn	01–09–2017 to 30–11–2017, 01–09–2018 to 30–11–2018, 01–09–2019 to 30–11–2019	10–09–2020 to 17–11–2020
Winter	01–01–2017 to 28–02–2017, 01–12–2017 to 28–02–2018, 01–12–2018 to 28–02–2019, 01–12–2017 to 31–12–2019	18–01–2020 to 25–01–2020



**Fig. 2.** One year's data (2020) for; a) sloped ERAD, Richborough; b) historical PV generation, Richborough; c) sloped ERAD, Grimsby; d) historical PV generation, Grimsby.

using modelled irradiance from Solar Radiation Data (SODA) as well as modelled weather data. As such, detailed information on the atmospheric conditions is required to determine irradiation. The critical limitation here is that atmospheric conditions are based on estimates which means that artificial errors may be introduced, thus decreasing the data's validity [18]. In addition, methods used to translate forecast solar irradiance values into PV power generation also introduces another layer of uncertainty into the prediction. This is due to numerous factors, such as the variation in PV installation, PV cell temperature and the sun's position at different times and days.

The proposed study aims to develop a new ML-based solar power forecasting model to forecast half-hourly PV generation. Unlike traditional approaches, the forecasting model development will not be based on the predicted local solar irradiance. Instead, accurately modelled sloped extra-terrestrial irradiance (ERAD) in conjunction with weather data will be used as input parameters. The parametric model for sloped ERAD does not depend on meteorological conditions. This approach aims to remove the predetermined relationship established in other parametric irradiation models. Instead, the model itself is allowed to find the relationship for the PV output.

The novelty and scientific significance of this paper is that it is the first reported study to use sloped ERAD as a key input parameter for PV output forecasting. The sloped ERAD based ML approach has been validated using historical data from two utility-scale solar PV plants operating under the UK's unpredictable weather conditions.

The advantage of this approach over other solar irradiance-based

approaches is that by using sloped ERAD the irradiance input is not weather dependent. As such, sloped ERAD can be accurately estimated at any time for any given PV installation, which is not possible for approaches based on modelled solar irradiance. To establish the effectiveness of using sloped ERAD quantitatively, a comparative study of the novel approach with the existing approach of using modelled solar irradiance was undertaken.

Another significance of the study is the modelling framework used, which includes data pre-processing, selection analysis of meteorological parameters for training and prediction, and model performance assessment. This offers a unique forecasting approach which can be applied in other locations under different weather conditions. Furthermore, the study paves the way for a simple and effective PV forecasting approach. This can be done through trials using the PV output at a larger scale in conjunction with forecast weather data to promote the appropriate use of this approach and highlight areas for future improvements.

The paper is structured as follow: Section 2 describes the case study, while the methodology is presented and explained in section 3. The results and overall performance are discussed in section 4, followed by conclusions drawn in section 5.

## 2. Sites, dataset, and visualisation

### 2.1. Site specifications and dataset

Four types of data were used for this study, all of which are

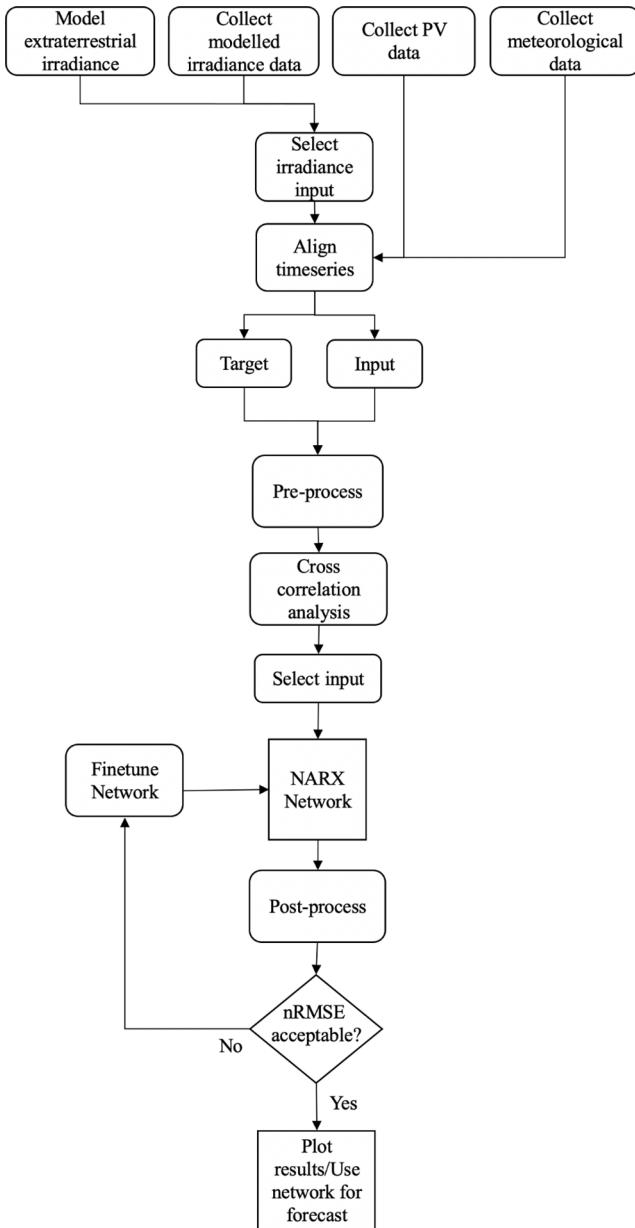


Fig. 3. Sequential procedure flowchart.

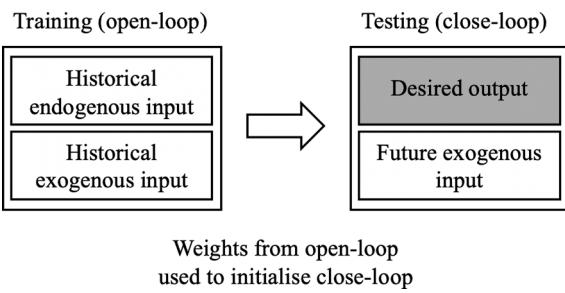


Fig. 4. Setup of NARXNN.

summarised in Table 1.

Data was collected for two locations with different characteristics: Richborough substation (Kent, England), Grimsby substation (Lincolnshire, England). Richborough has a favourable position in the UK for solar PV readings and currently has a capacity of 102.9 MWp. Grimsby

**Table 3**  
Endogenous and exogenous inputs.

Exogenous Inputs	Endogenous Inputs	Time Lag
Sloped ERAD Cloud-cover Visibility Temperature	Historical PV output	Seven days

has a capacity of 53 MWp with weather of a more intermittent nature.

Fig. 1 show the locations of the two solar plants and weather stations in the UK.

## 2.2. Training and validation dataset

Data for four years (1st of January 2017 to 31st of December 2020) was collected at the two locations. Data for testing and training was extracted as shown in Table 2, creating sixteen datasets: eight for each site. The dataset is divided according to literature [9].

## 2.3. Data visualisation

Fig. 2 depicts one year (2020) of the dataset for Richborough (Fig. 2a-b) and Grimsby (Fig. 2 c-d). This visualisation of data allows for studying changes throughout the day. The data illustrates that irradiance and PV peak in summer and dip in winter. Sloped ERAD is close to identical for the two sites due to their geographical proximity. The impact of the atmosphere on solar irradiance is reflected in the PV output. Furthermore, note how the daily time range for output is shorter during winter than during summer. This is due to solar geometry; as the position of the sun changes, so does the maximum solar radiation available.

## 3. Methodology

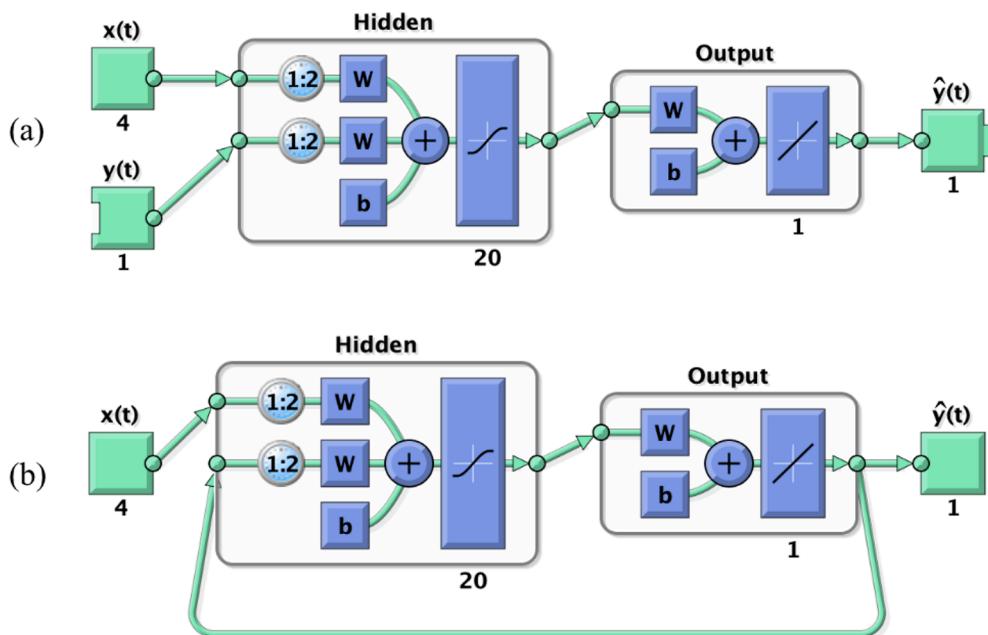
Fig. 3 illustrates the applied sequential model approach used for forecasting PV output power. Night-values were removed by eliminating Zeros from the training and validation datasets, which were normalised between 0 and 1 to improve precision, reduce regression error, and sustain correlation. All neural network simulations are performed using MATLAB R2020b. This section describes the methodology based on *sloped ERAD* being used as the selected irradiance input. For the comparative approach, the modelled irradiance in Table 1 was instead selected as the irradiance input in Fig. 3.

NARXNN was selected as the model based on findings by Su, Batzelis & Pal [9]. Results suggested that the NARXNN was the best performing model out of ten conventional ML-based models, which were all compared on the same framework. The NARXNN is widely used globally for solar PV predictions, thanks to its superiority to other models. For example, the NARXNN was used for PV output forecasting by Buritrago & Asfour [22] in the United States and by Vaz, Elsinga, van Sark & Brito [23] in the Netherlands.

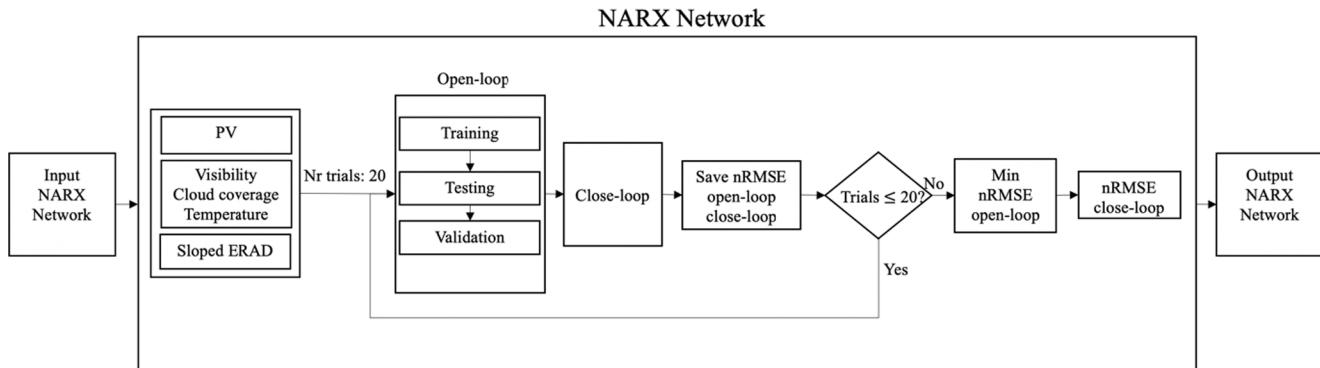
Prior to using the NARXNN model for PV output forecasting, the following steps were completed:

- Data pre-processing, including the alignment of data time-series.
- Defining model target and selection of corresponding data for training and prediction.
- Correlation analysis of PV output against *sloped ERAD* and meteorological data (*cloud-cover*, *visibility*, *temperature*, *wind speed* and *wind direction*).

The correlation analysis ensures that the high-impact meteorological parameters are selected as input parameters for PV forecast, as these parameters might be site-specific. Data was evaluated using Pearson's correlation coefficient during the correlation analysis, as shown in



**Fig. 5.** NARXNN training models with 20 hidden neurons, 1:2 input delays, and 1:2 feedback delays; a) open-loop; and b) close-loop, where w stands for weights and b stands for bias.



**Fig. 6.** Proposed network construction.

Equation (1).

$$r = \frac{\text{cov}(X, Y)}{\sqrt{\text{cov}(X, X)} \sqrt{\text{cov}(Y, Y)}} \quad (1)$$

where  $r > 0$  means two parameters have a correlation and are directly related to one another, and  $r < 0$  means two parameters are inversely related. When  $r$  is close to 1, both have a close relationship, and when  $r = 0$ , no correlation is found. The highly correlated parameters to PV output were then selected for the NARXNN.

### 3.1. Network overview

Two different architectures can be employed with the NARXNN model, where the first performs *one-step-ahead* prediction and the second performs *multi-step-ahead* prediction. The former is often referred to as *open-loop* and the latter is referred to as *close-loop* [9]. The configurations may also be referred to as series-parallel architecture and parallel architecture, respectively [24].

The model uses endogenous and exogenous inputs, which are the internal and external inputs, respectively. In this case, PV output was treated as endogenous input, whereas *sloped ERAD* and weather data were treated as exogenous inputs. Fig. 4 shows a simplified setup of the

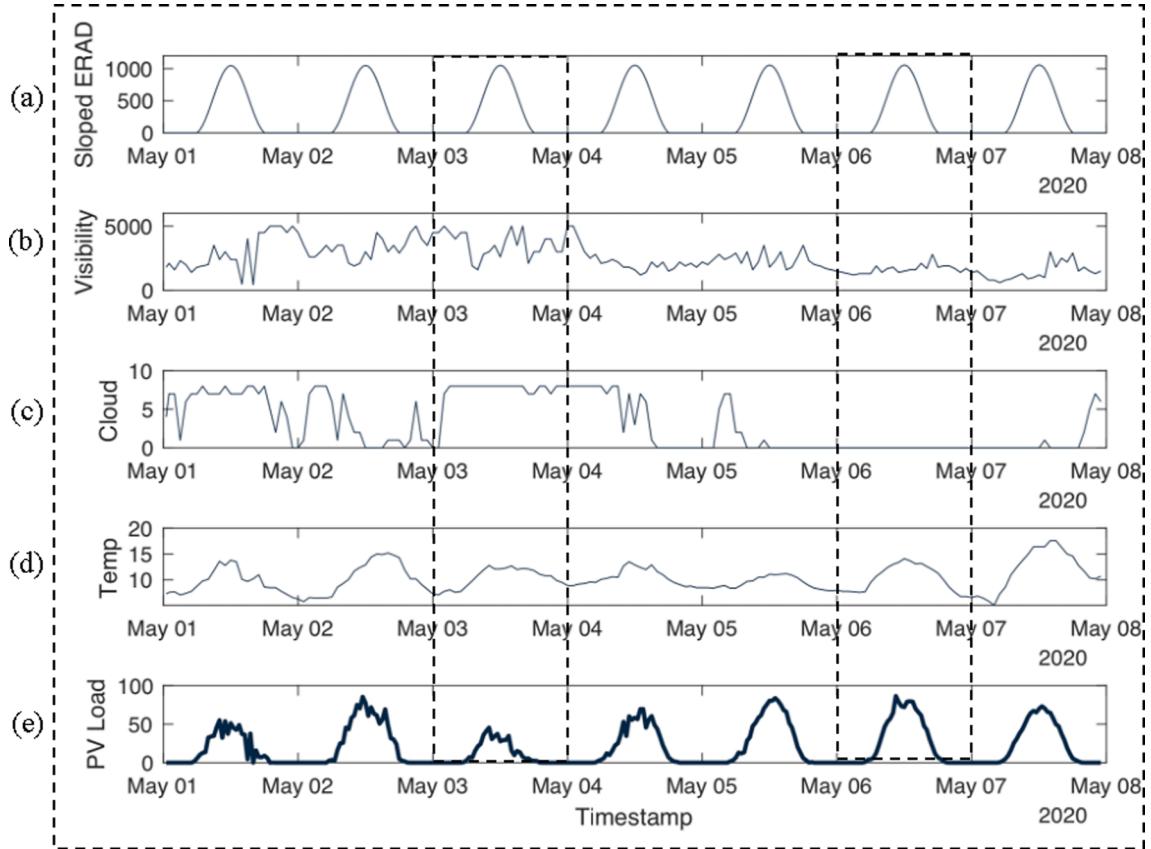
NARXNN where both architectures are used.

The selected meteorological parameters served as exogenous inputs were based on the correlation analysis between a set of weather data and PV output. The results are presented in the next section. Table 3 details the endogenous and exogenous inputs for the NARXNN model.

Ideally, forecast data should be used for the meteorological parameters used when forming future exogenous inputs. Due to the data's availability, historical meteorological data was split into two parts in this work, with data prior to 2020 treated as historical inputs and data from 2020 treated as future exogenous inputs. This selection also allows for testing the model's performance when using the same approach as taken in the literature [14, 9].

The models for the NARXNN configurations used in this study are presented in Fig. 5. The NARXNN architecture for both configurations consisted of one endogenous input node  $y(t)$  four exogenous input nodes  $x(t)$ , one hidden layer with 20 nodes and one output node. The endogenous input node  $y(t)$  for *open-loop* is set up in the same way as for  $x(t)$ , meaning that a new  $y(t)$  is needed for each point from the historical dataset (Fig. 5a). However, for *close-loop* the value obtained as the predicted endogenous output  $\hat{y}(t)$  was instead fed back for the next prediction (Fig. 5b).

The commonly applied workflow, which was adopted, is to create the



**Fig. 7.** Photovoltaic power characteristics and meteorological factors for Richborough; a) sloped ERAD ( $\text{W}/\text{m}^2$ ); b) visibility (m); c) cloud-cover (okta); d) temperature ( $^\circ\text{C}$ ); e) PV load (MW).

network in *open-loop* (Fig. 5a). The *open-loop* uses *one-step-ahead* predictions, making it suitable for training the model. The weights and biases from the *open-loop* were then used to initialise *close-loop* predictions [25]. In Fig. 5b the four exogenous input nodes consist of future weather and *sloped ERAD* inputs. Together these streams of input created the first predicted PV output. This value was fed back and used for the next point, to allow for *multi-step-ahead* predictions. Both input delays and feedback delays were set to 1:2. Full interconnection was selected. The sigmoid function was used as the activation function, and the Levenberg-Marquart backpropagation was used as the learning algorithm in accordance with literature [9].

### 3.2. Network construction

Fig. 6 illustrates a more detailed flowchart for the proposed network construction. Here, the pre-processed inputs were fed into the network, where the dataset is split into training, testing and validation (in a 70%, 15% and 15% split, respectively) in accordance with literature [22]. This split is the default setting and allows for a large amount of data being used for training. The training set used Levenberg-Marquardt backpropagation in an *open-loop* configuration to determine the weights of the training set. The loop was closed and initialised with the weights from the *open-loop*. The model then used the created PV output prediction as an input to generate the next forecast value.

The normalised root mean square error (nRMSE) was employed as an evaluation method, as shown in equation (2). It has been widely used in literature for substations [26]. For example, [27] uses nRMSE to compare results for different models for PV output forecasting at a power station. By introducing the normalised prediction error, with the installed capacity  $P_{ins}$ , comparisons between errors referring to different solar farms or substations can be made.

$$\text{nRMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N \left( \frac{\hat{P}(i) - P(i)}{P_{ins}} \right)^2} \quad (2)$$

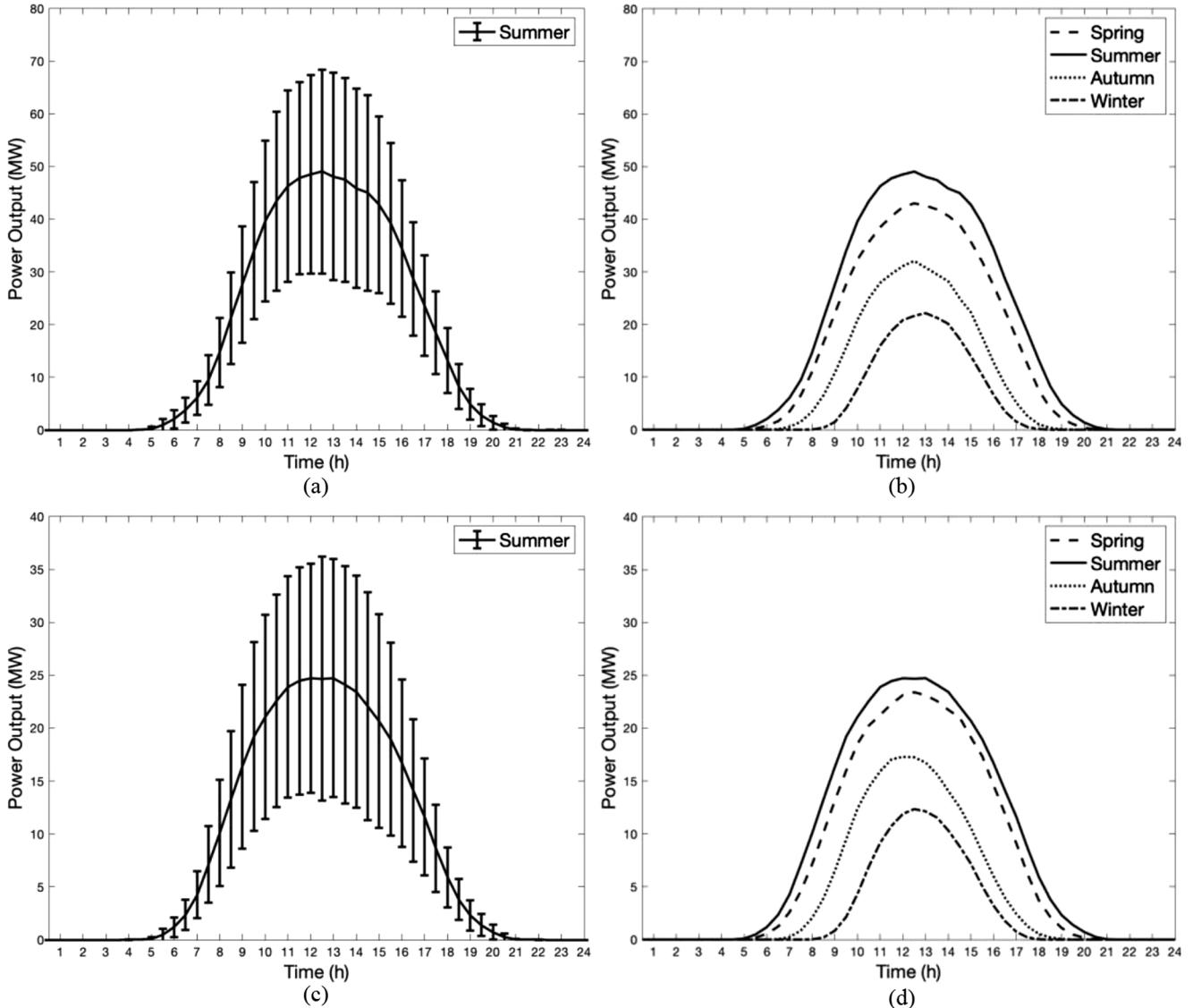
where  $N$  is the number of data samples;  $\hat{P}(i)$  and  $P(i)$  represent the predicted and measured power at the time  $i$ .  $P_{ins}$  is the installed capacity.

nRMSE was saved for both configurations, and the process was repeated 20 times (which was selected as an arbitrary number). The network's forecast (*close-loop*) performance was determined based on the lowest training (*open-loop*) nRMSE. The process was repeated for the comparative approach, where *sloped ERAD* was changed to *modelled irradiance*.

## 4. Data analysis and forecast performance

### 4.1. Meteorological analysis and PV output

Fig. 7 shows *sloped ERAD*, *visibility*, *cloud-cover*, *temperature*, and PV output over a week in May at Richborough. The *sloped ERAD* exhibits a perfect bell curve at the top of the atmosphere due to being unaffected by atmospheric weather conditions (Fig. 7a). The radiation striking a surface on earth can be expressed as the residual of *sloped ERAD* once it has passed through the atmosphere, which strongly correlates to the PV output. Meteorological variables (Fig. 7b-d) clearly demonstrate that weather conditions influence the amount of radiation that reaches the surface of the PV panel. It is evident that low PV output occurs on cloudy days (i.e., on the 3rd of May) with a weak correlation to *sloped ERAD*, while high PV output is observed on clear sky days (i.e., on the 6th of May) with a strong correlation to *sloped ERAD*.



**Fig. 8.** Seasonal data over three years for; a) Richborough, mean summer PV output with standard deviation; b) Richborough, mean seasonal PV output for all seasons; c) Grimsby, mean summer PV output with standard deviation; and d) Grimsby, mean seasonal PV output for all seasons.

#### 4.2. Seasonal analysis

An analysis of PV output for the two locations shows that the PV output varies significantly even within the summer season. The standard deviations are up to 50% of the peak hourly mean outputs (Fig. 8a & 8c). This suggests that unpredictable weather conditions significantly impact PV output, reinforcing the importance of developing a PV forecasting tool to balance supply and demand. The seasonal data shows the seasonal variation of the mean PV output for both sites, with the highest output in summer followed by spring, autumn and winter (Fig. 8b & 8d). However, the variance between summer and spring is not significant as peak outputs for both sites are at a range of 43–50% of the installed capacity throughout these seasons (Richborough: 103 MWp and Grimsby: 53 MWp).

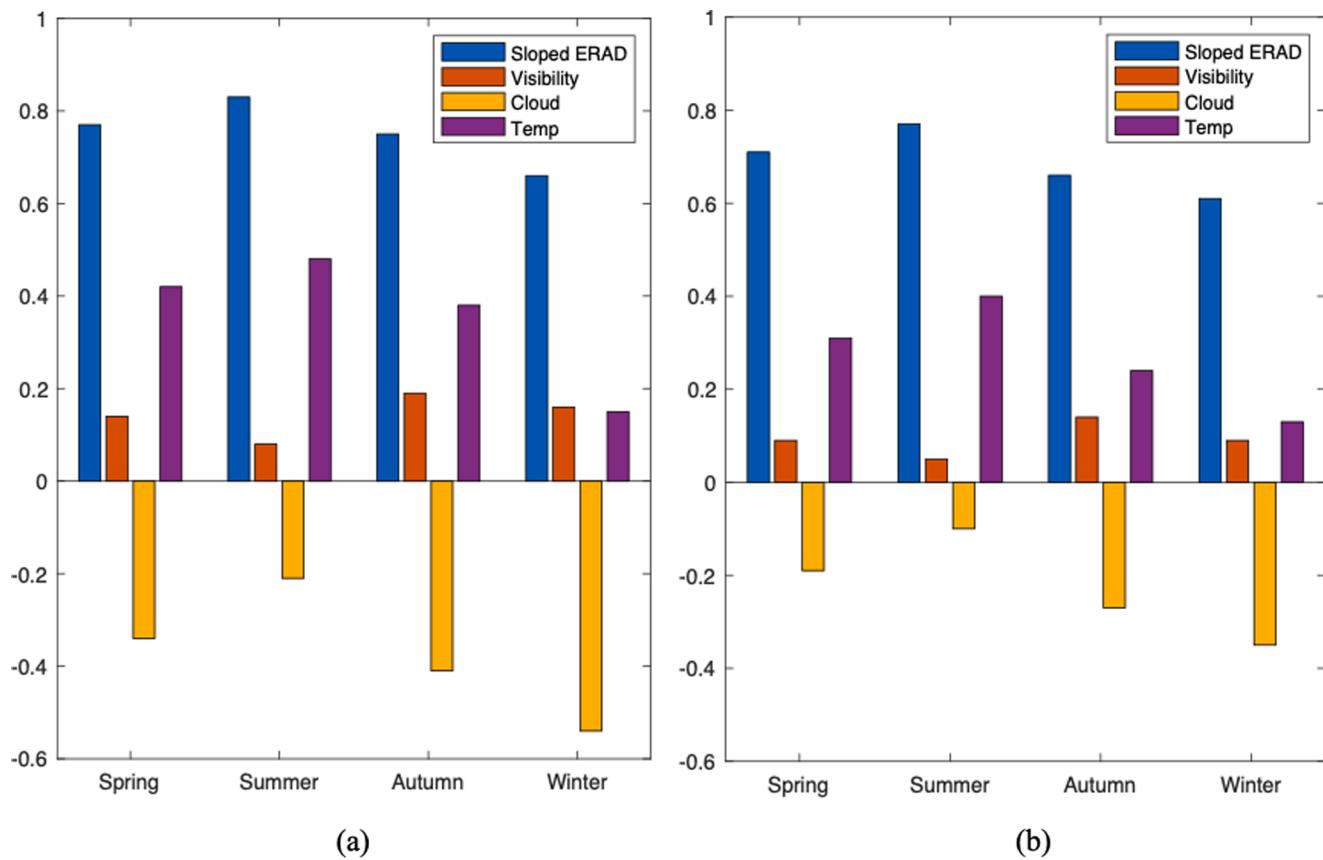
#### 4.3. Correlation coefficient analysis

The correlation analysis shows that Pearson's correlation coefficient varies with both season and location. In general, the highest correlation is found for *Sloped ERAD*, followed by *cloud-cover*, *temperature*, *visibility*, *wind speed* and *wind direction*. Fig. 9 presents the correlation coefficients of the four variables that correlate most closely to PV output: *Sloped ERAD*, *cloud-cover*, *temperature* and *visibility*. A comparison of the two

locations suggests that all selected variables have a stronger correlation to Richborough (Fig. 9a) than to Grimsby (Fig. 9b).

As expected, *sloped ERAD* has the highest positive correlation with PV output for all seasons. This justifies the use of *sloped ERAD* as an essential input parameter for the model. *Cloud-cover* shows a high negative correlation to PV output, although it varies significantly with the seasons and PV locations. This is because *cloud-cover* reduces the amount of solar irradiance reaching a PV surface, and its impact depends on the PV local weather conditions, as revealed by numerous studies [17,10].

It is interesting to note that temperature shows a positive correlation to PV output. As the PV's efficiency decreases with the increased temperature, a plausible explanation for this positive correlation is that a high temperature also generally correlates to high solar irradiance on PV given the UK's temperate weather conditions. As such, *temperature* is related to *sloped ERAD* and is not an independent variable. The highest temperature correlation to PV output occurs in summer, followed by spring, autumn and winter (Fig. 9). *Visibility* has the weakest correlation to the PV output among these four parameters for all seasons, regardless of location. The correlation analysis provides information for the model parameter selection and rationale for the model sensitivity analysis, presented in the later section.



**Fig. 9.** Pearson's correlation coefficient between PV and exogenous input variables for a) Richborough; and b) Grimsby.

**Table 4**  
Seasonal forecast performance for Richborough and Grimsby (nRMSE).

Season	Richborough		Grimsby	
	Irrad approach <sup>1</sup>	ERAD approach <sup>2</sup>	Irrad approach	ERAD approach
Spring	5.19%	6.36%	8.96%	8.95%
Summer	5.95%	5.54%	5.30%	8.26%
Autumn	6.87%	6.46%	6.53%	9.35%
Winter	3.85%	5.25%	4.26%	5.70%

1. Irrad approach uses modelled irradiance, visibility, cloud-cover and temperature as input parameters.

2. ERAD approach uses sloped ERAD, visibility, cloud-cover and

#### 4.4. Forecast performance

The forecast performance (*close-loop*) for the two approaches is summarised in [Table 4](#), where the nRMSE for a 7-day-ahead forecast for Richborough and Grimsby is presented. Both approaches use *visibility*, *cloud-cover* and *temperature* as meteorological input parameters. The *Irrad approach* is based on modelled irradiance, which was used as a comparative study for the sloped *ERAD based approach*.

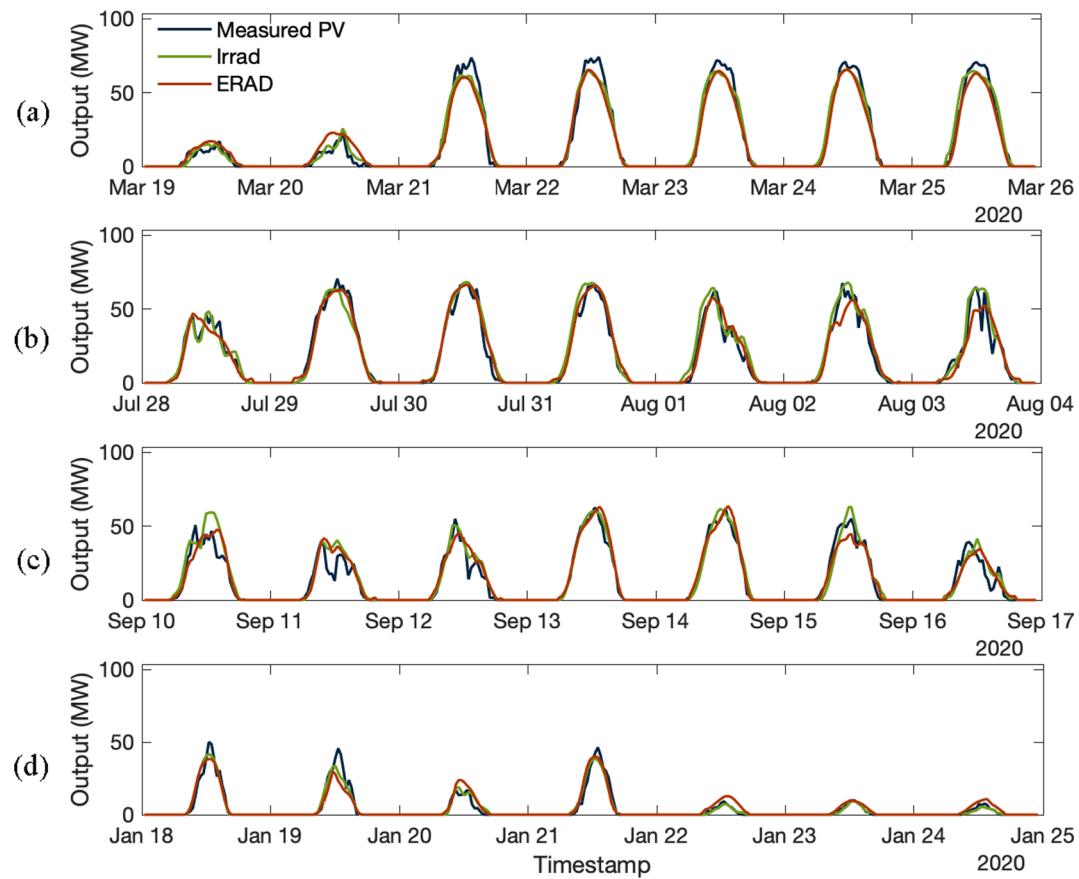
There are trends that generally apply to both approaches. Forecasting in summer yields better results than in spring and autumn for both locations, which was expected due to the steady weather and a large amount of solar irradiance. This is very useful since the peak PV output is generated during this season. Forecasts in spring and autumn show very similar results; the two locations do not show concordance in relation to which of the seasons is the most difficult to forecast. Therefore, in this case, the complexity of the daily pattern will have a greater impact on forecasting than the season itself.

During the winter, shorter sunshine periods mean that data is

limited. It was therefore surprising that winter achieved the lowest nRMSE for both locations. However, care must be taken when analysing nRMSE, as the evaluation is based on capacity. Winter is the season where the measured PV output is the lowest, and the contrast between the measured output and the installed capacity will therefore be the largest. For this reason, the nRMSE for low PV outputs (such as those experienced during the winter) does not necessarily indicate high accuracy for the PV output forecast, and a graphical representation is therefore needed to confirm the results. The measured PV output for three of the winter days was below 10% of the total capacity, which is the most likely reason for the good results. Furthermore, the strongest correlation with *cloud-cover* for both locations was for winter, which may further explain the results.

A seasonal representation of the measured and forecasted PV output results are presented in [Figs. 10 and 11](#) (for Richborough and Grimsby, respectively). Results show an excellent PV output forecast for all four seasons at Richborough ([Fig. 10](#)). In particular, the forecast captured huge intra-daily variations of PV output. For example, the peak PV output on the 19th and 20th of March 2020 (~20 MW) is less than one-third of the following five-day outputs (~60 MW) in spring ([Fig. 10a](#)). Moreover, the daily peak output variation during the winter week fluctuates (~40 MW) ([Fig. 10d](#)). For Grimsby, the forecast is less accurate when compared with Richborough. However, the forecast can capture the daily variation of PV output reasonably well. This is evident in the forecasting performance on the 19th and 20th of March in spring, 15th and 16th of September in autumn and 21st and 24th of January in winter ([Fig. 11](#)). This forecasting feature is particularly useful to balance the supply and demand of the electricity system.

The *ERAD approach* does not perform well in capturing the inter-hourly variability of PV output due to the rapid change to weather conditions. In fact, capturing the inter-hourly variability of PV output is very challenging given the UK's erratic weather. A plausible explanation for this is the



**Fig. 10.** Forecast performance for Richborough; a) spring; b) summer; c) autumn; d) winter.

smoothing of weather data, since linear interpolation was applied for the meteorological variables to create a *half-hourly* dataset from an *hourly* dataset. Using smoothed weather data as input parameters for modelling, training, and forecasting fails to accurately represent the local weather conditions reflected by the measured PV output. However, it should be recognised that even in situations where high-frequency weather data is available, the data must represent the local weather conditions for effective PV output modelling. High-spatial-resolution weather data is often not available in the UK, given the island's weather conditions. Overall, this study demonstrates that the *ERAD approach* can capture a clear daily trend, which ultimately is more valuable for stakeholders.

When comparing the *Irrad approach* with the *ERAD approach*, there are a number of things to take into consideration. Quantitatively the *Irrad approach* is more effective in capturing inter-hourly fluctuations compared to the *ERAD approach*. This is because the modelled irradiance was obtained based on high frequency (15 min) historical weather data through satellite images. Measured weather data in the UK is readily available, but in lower frequency intervals (hourly), and the *Sloped ERAD* is dependent on this information. It is therefore not surprising that the *Irrad approach* outperforms the *ERAD approach* under these circumstances.

The key limitation when using the *Irrad approach* is that it requires tremendous resources for accurate weather forecasting at 15-minute frequency intervals (to generate *modelled irradiance*) in addition to requiring values for *visibility*, *cloud-cover* and *temperature*, forecast at hourly frequency intervals for the coming week ahead. *Sloped ERAD* on the other hand is unaffected by weather conditions and can be determined with high accuracy, both for past and future values. The *ERAD approach* only requires hourly forecasting of *visibility*, *cloud-cover* and *temperature*, making it suitable to use in practice.

Richborough exhibits a generally satisfying performance for all seasons compared to Grimsby, where the performance is less accurate. This suggests that forecast accuracy does depend on local meteorological

conditions and plant capacity. As *sloped ERAD* was close to identical for both substations, due to their proximity, meteorological conditions will inevitably play an integral role in subsequent predictions using *sloped ERAD*. When observing the measured PV output for the two substations, fluctuations were more common for Grimsby, with frequent midday spikes. The measured data revealed that the weather conditions at Grimsby are more intermittent than Richborough due to its location in the UK. Unsurprisingly, the frequent variations at Grimsby are much more challenging to forecast.

In addition to the variability in weather, variation in capacity may also be an important factor. As suggested from the results, the approach seems to be more effective for a larger PV installation. This may be because microclimate has a more significant impact on smaller PV installations such as Grimsby. Richborough, on the other hand, is less affected by unexpected cloud-cover or weather changes. This information is useful because accurate 7-day-ahead PV forecasting is needed to minimise curtailment on a regional or national level. Furthermore, this incentivises an increase in connection of solar PV panels to substations; providing better forecast predictions in addition to an uptick in electricity production. This is valuable information for stakeholders in the planning of future PV installations.

#### 4.5. Sensitivity analysis of input parameters

A preliminary sensitivity analysis was conducted using several combinations of input parameters in order to validate the selection of the meteorological input parameters for the *ERAD approach*. This was done to optimise the number of input parameters for the model. These combinations were constructed based on the correlation analysis (Fig. 9). *Sloped ERAD* is included for all combinations as it has the highest correlation to PV output when compared all other parameters as well as superior availability (it can be modelled to high accuracy).

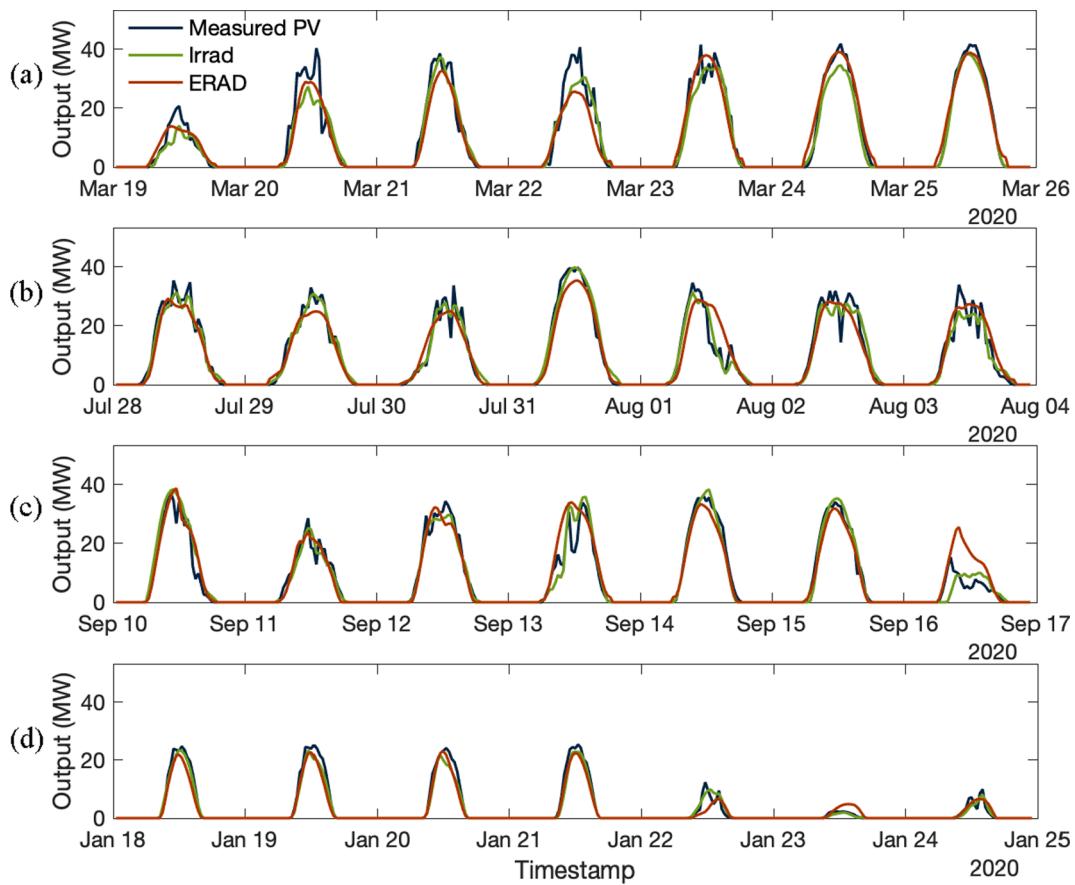


Fig. 11. Forecast performance for Grimsby; a) spring; b) summer; c) autumn; d) winter.

Table 5

Forecast performance using different combinations of input parameters for summer (nRMSE).

Combination	Exogenous Input	Richborough	Grimsby
A	Sloped ERAD + Cloud-cover + Temperature + Visibility	5.65%	8.85%
B	Sloped ERAD + Cloud-cover + Temperature	5.59%	8.58%
C	Sloped ERAD + Cloud-cover	5.68%	8.25%
D	Sloped ERAD + Temperature	6.83%	11.33%
E	Sloped ERAD + Cloud-cover + Temperature + Wind direction	5.64%	8.62%
F	Sloped ERAD + Cloud-cover + Temperature + Wind speed	5.27%	8.78%

regardless of location). Each combination's forecast performance was assessed using summer data, as the peak PV output is found in this season.

Results show that Combination C, which used *sloped ERAD* and *cloud-cover* only as input parameters, is the optimal option for PV forecasting (Table 5). This is because the model performance was either maintained or increased for both locations when temperature and visibility were removed from input parameters, resulting in a simpler and more cost-effective model. The sensitivity analysis confirms that *sloped ERAD* and *cloud-cover* are the most important input parameters. The impact of other meteorological variables (e.g., *temperature*, *visibility*, *wind speed* and *wind direction*) on the model performance are negligible, as indicated by the correlation analysis. This reinforces that the correlation analysis for the input parameter selection is critical for a successful forecast, as the model input parameters can be site-specific.

#### 4.6. Model development and potential applications

Ultimately, the aim was to create a model that will work in practice and that can pave the way for solving industry issues related to the absence of recorded solar irradiance globally. The use of modelled *sloped ERAD* to forecast PV output means that the novel approach removes complexity and potential errors related to other modelled irradiance parameters whilst still maintaining a high accuracy. Adopting *sloped ERAD* as an input parameter in future research means that reliable PV forecasting can be obtained, which is beneficial for future development in this field.

The novel approach will require further testing using weather forecasting data. The model will not perform as effectively as suggested in Figs. 10 and 11 if using forecast weather data in place of historical weather data. However, to make the model practical the use of forecasted weather data is essential. The forecast horizon will inevitably impact the accuracy of the weather forecast. For example, for a five-day and a seven-day forecast, 90% and 80% accuracy can be expected, respectively. For a 10-day (or longer) forecast, equivalent accuracy values are only obtained half of the time [28]. 7-day-ahead weather forecast data is available from the Met Office. The granularity of cloud data that is available to the public is less accurate than that used for recorded data. However, the forecast of *cloud-cover* in oktas, similar to historical data available from the Met Office, is available as a fee-paid service [29].

The model would be particularly useful for larger scale applications (i.e., county or regional scale). This is because the intra-hourly change of PV output caused by the rapid change in local weather conditions can be smoothed. Potentially a more accurate intra-hourly PV forecast can be achieved on this scale. Trials using the PV output at a larger scale in conjunction with forecast weather data promote the appropriate use of

this approach and highlight areas for future improvements.

## 5. Conclusion

In the present paper, a novel machine learning approach for solar photovoltaic energy output forecasting was developed using *sloped extra-terrestrial irradiance*. The main conclusions are detailed below.

- The promising results validate the concept of using *sloped extra-terrestrial irradiance* as an input parameter in conjunction with meteorological variables. The uniqueness of this approach relates to the approach's ability to precisely model *sloped extra-terrestrial irradiance* – an input parameter that has a stronger correlation to PV output compared to widely recorded meteorological parameters. Furthermore, since *sloped extra-terrestrial irradiance* is based on solar geometry and therefore not influenced by atmospheric weather conditions, it can also be accurately modelled for future inputs. This is a great advantage compared to conventional approaches where inaccurately predicted solar irradiance is often used.
- The novel forecasting model captures huge intra-daily variations in photovoltaic output, which is particularly useful for balancing supply and demand of the electricity system.
- The use of *sloped extra-terrestrial irradiance* and *cloud-cover* presented the optimal combination of input parameters as these provided the simplest and most cost-effective model without reducing accuracy.
- The model was tested and validated in what can be considered a worst-case scenario in UK's unpredictable and intermittent weather conditions. The excellent model performance suggests that there is now a strong imperative to use the model in other locations where weather conditions are more stable. The proposed approach offers universal value as it only requires coordinates and weather data.

Trials using photovoltaic output at a larger scale in conjunction with forecast weather data promote the appropriate use of this approach and highlight areas for future improvement, including consideration of *wear and tear* as an input parameter to represent reduced PV efficiency over time.

## CRediT authorship contribution statement

**Cornelia A. Fjelkestam Frederiksen:** Data curation, Methodology, Formal analysis, Software, Validation, Visualisation and Writing-original draft. **Zuansi Cai:** Conceptualization, Data curation, Funding acquisition Methodology, 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.

## Acknowledgement

The authors wish to thank the editors and reviewers for their valuable feedback. Their constructive comments have improved this paper. This work was partly supported by Edinburgh Napier University through Internal Funding Competition.

## References

- [1] Committee on Climate Change. Net Zero The UK's contribution to stopping global warming. Committee on Climate Change; 2019.
- [2] Dobson-Smith J. 175MW of new PV deployed in first quarter of 2021, Solar Energy UK; 15 April 2021. [Online]. Available: <https://solarenergyuk.org/news/175mw-of-new-pv-deployed-in-first-quarter-of-2021/> [accessed 05 July 2021].
- [3] Cockburn H. UK ranked sixth in world for share of electricity generated by wind and solar. Independent; 8 July 2021. [Online]. Available: <https://www.independent.co.uk/climate-change/news/renewable-energy-rank-world-list-b1880639.html> [accessed 16 July 2021].
- [4] van der Meer DW, Shepero M, Svensson A, Widén J, Munkhammar J. Probabilistic forecasting of electricity consumption, photovoltaic power generation and net demand of an individual building using Gaussian Processes. Appl Energy 2018;213: 195–207. <https://doi.org/10.1016/j.apenergy.2017.12.104>.
- [5] Winterman D. Is the British weather unique in the world?; 2013. [Online]. Available: <https://www.bbc.com/news/magazine-24305230> [accessed 5 March 2021].
- [6] National Grid. "Summer Outlook Report 2018. National Grid; 2018.
- [7] Sheffield Solar. National PV Generation; 2021. [Online]. Available: <https://www.solar.sheffield.ac.uk/pvlive/> [accessed 2 April 2021].
- [8] Renewable Energy Foundation. Balancing mechanism wind farm constraint payments; 2020.
- [9] Su D, Batzelis E, Pal B. Machine learning algorithms in forecasting of photovoltaic power generation. International Conference on Smart Energy Systems and Technologies (SEST) 2019:1–6. <https://doi.org/10.1109/SEST.2019.8849106>.
- [10] Ahmed R, Sreeram V, Mishra Y, Arif MD. A review and evaluation of the state-of-the-art in PV solar power forecasting: Techniques and optimization. Renew Sustain Energy Rev 2020;(124). <https://doi.org/10.1016/j.rser.2020.109792>.
- [11] Nespoli A, Ogliari E, Leva S, Massi Pavan A, Mellit A, Lugh Vi, et al. Day-Ahead Photovoltaic Forecasting: A Comparison of the Most Effective Techniques. Energies 2019;12(9):1621. <https://doi.org/10.3390/en12091621>.
- [12] Barberi F, Rajakaruna S, Ghosh A. Very short-term photovoltaic power forecasting with cloud modelling: A review. Renew Sustain Energy Rev 2017;75:242–63. <https://doi.org/10.1016/j.rser.2016.10.068>.
- [13] Li P, Zhou K, Lu X, Yang S. A hybrid deep learning model for short-term PV power forecasting. Appl Energy 2020;(259). <https://doi.org/10.1016/j.apenergy.2019.114216>.
- [14] Andersson WW, Yakimenko OA. Using neural networks to model and forecast solar PV power generation at Isle of Eigg. CPE-POWERENG; 2018. p. 1–8. <https://doi.org/10.1109/CPE.2018.8372522>.
- [15] Liu L, Liu D, Sun Q, Li H, Wennersten R. Forecasting power output of photovoltaic system using a BP network method. Energy Procedia 2017. <https://doi.org/10.1016/j.egypro.2017.12.126>.
- [16] Das UK, Tey KS, Seyyedmahmoudian M, Melkilef S, Idris MYI, Van Deventer W, et al. Forecasting of photovoltaic power generation and model optimization: A review. Renew Sustain Energy Rev 2018;81:912–28. <https://doi.org/10.1016/j.rser.2017.08.017>.
- [17] Munee T, Gueymard C, Kambezidis H. Solar Radiation and Daylight Models. 2nd ed. Oxford: Taylor & Francis Group; 2004.
- [18] Ahmad A, Anderson T. Global solar radiation prediction using artificial neural network models for New Zealand;Solar2014: The 52nd Annual Conference of the Australian Solar Council.
- [19] Sheffield Solar. Regional PV Generation; 2021. [Online]. Available: <https://www.solar.sheffield.ac.uk/pvlive/regional/> [accessed 03 January 2021].
- [20] CEDA Archive. Dataset MIDAS: UK Hourly Weather Observation Data; 2021. [Online]. Available: <https://catalogue.ceda.ac.uk/uuid/916ac4bbc46f7685ae9a5e10451bae7c?jump=related-docs-anchor> [accessed 1 April 2021].
- [21] CAMS. CAMS solar radiation time-series; 2021. [Online]. Available: <https://ads.atmosphere.copernicus.eu/cdsapp#!/dataset/cams-solar-radiation-timeseries?tab=overview> [accessed 24 September 2021].
- [22] Buitrago J, Asfour S. Short-Term Forecasting of Electric Loads Using Nonlinear Autoregressive Artificial Neural Networks with Exogenous Vector Inputs. Energies 2017;10(40). <https://doi.org/10.3390/en10010040>.
- [23] Vaz AGR, Elsinga B, van Sark WGJHM, Brito MC. An artificial neural network to assess the impact of neighbouring photovoltaic systems in power forecasting in Utrecht, the Netherlands. Renew Energy 2016;85:631–41. <https://doi.org/10.1016/j.renene.2015.06.061>.
- [24] Boussaada Z, Curea O, Remaci A, Camblong H, Bellaaj NM. A Nonlinear Autoregressive Exogenous (NARX) Neural Network Model for the Prediction of the Daily Direct Solar Radiation. Energies 2018. <https://doi.org/10.3390/en11030620>.
- [25] MathWorks. Design Time Series NARX Feedback Neural Networks. MathWorks; 2021. [Online]. Available: <https://www.mathworks.com/help/deeplearning/ug/design-time-series-narx-feedback-neural-networks.html> [Accessed 19 July 2021].
- [26] Madsen H, Pinson P, Kariniotakis G, Nielsen HA, Nielsen TS. Standardizing the Performance Evaluation of Short-Term Wind Power Prediction Models. Wind Eng 2005;29(6):475–89. <https://doi.org/10.1260/030952405776234599>.
- [27] Huang Y, Lu J, Liu C, Xu X, Wang W, Zhou X. Comparative study of power forecasting methods for PV stations. 2010 International Conference on Power System Technology 2010:1–6. <https://doi.org/10.1109/POWERCON.2010.5666688>.
- [28] NOAA. How Reliable Are Weather Forecasts?; 2021. [Online]. Available: <https://scijinks.gov/forecast-reliability/> [accessed 4 March 2021].
- [29] Met Office. Personal communication; 2021.