

Article

***k*-Nearest Neighbor Neural Network Models for Very Short-Term Global Solar Irradiance Forecasting Based on Meteorological Data**

Chao-Rong Chen and Unit Three Kartini *

Department of Electrical Engineering, National Taipei University of Technology, 1, Section 3, Zhong-Xiao (Chung-Hsiao) E. Rd., Da'an Dist., Taipei 106, Taiwan; crchen@ntut.edu.tw

* Correspondence: uunitthree@gmail.com or t101319021@ntut.edu.tw; Tel.: +886-2-27712171 (ext. 2112)

Academic Editor: Silvio Simani

Received: 24 November 2016; Accepted: 1 February 2017; Published: 8 February 2017

Abstract: This paper proposes a novel methodology for very short term forecasting of hourly global solar irradiance (*GSI*). The proposed methodology is based on meteorology data, especially for optimizing the operation of power generating electricity from photovoltaic (PV) energy. This methodology is a combination of *k*-nearest neighbor (*k*-NN) algorithm modelling and artificial neural network (ANN) model. The *k*-NN-ANN method is designed to forecast *GSI* for 60 min ahead based on meteorology data for the target PV station which position is surrounded by eight other adjacent PV stations. The novelty of this method is taking into account the meteorology data. A set of *GSI* measurement samples was available from the PV station in Taiwan which is used as test data. The first method implements *k*-NN as a preprocessing technique prior to ANN method. The error statistical indicators of *k*-NN-ANN model the mean absolute bias error (*MABE*) is 42 W/m² and the root-mean-square error (*RMSE*) is 242 W/m². The models forecasts are then compared to measured data and simulation results indicate that the *k*-NN-ANN-based model presented in this research can calculate hourly *GSI* with satisfactory accuracy.

Keywords: global solar irradiance (*GSI*); photovoltaic (PV); very short term; forecasting; *k*-nearest neighbor (*k*-NN); artificial neural network (ANN)

1. Introduction

Nowadays, forecasting global solar irradiance (*GSI*) is an essential task, particularly related to the increased use of photovoltaic (PV) solar energy as a power source. Forecasting *GSI* can be executed in different terms: long-term, medium-term, and short-term. Since solar power is categorized as an intermittent energy source, forecasting is paramount to regulate electricity loads in power networks. It also functions to optimize power delivery and unit commitment and by extension, it helps minimize the operating costs of power systems [1]. With a forecast, it is expected that plant operation control systems can be improved, so as to balance power generation and load. Moreover, distribution of load, electric energy storage, and energy supply will be maximized and more reliable.

The performance of PV systems is heavily influenced by meteorological conditions such as temperature, global irradiation, humidity, wind speed and wind direction [2]. The relation is clear: electrical energy generated by the PV solar depends on the amount of the *GSI* received by the PV panels. Solar irradiance absorbed in each PV panel varies depending on geographic location, time, and the absorption capacity of the PV panels. Previous studies have presented a variety of mathematical models for *GSI* forecasting in relation to meteorological variables. *GSI* forecasting with *k*-nearest neighbor (*k*-NN) statistical methods has been described [3,4]. Hocaoglu [3] and Pedro and Coimbra [4] presented modeling of solar irradiation with stochastic methods using a *k*-NN artificial neural network (ANN) at a PV station. Solar irradiation prediction is an important problem in

geosciences, with direct applications in renewable energy, where the data in the form of time series can also be analyzed using a regression model as described in reference [5]. Salcedo-Sanz et al. [5] worked on the prediction of daily global irradiation using a temporal Gaussian process, in which the study explains the suitability of Gaussian regression (GPR) for the estimation of solar irradiation compared to other machine learning regression algorithms. The *GSI* forecasting is not only used in stochastic modeling, but in other studies [6,7] attempted forecasting was analyzed using exponential smoothing combined with decomposition methods and least absolute shrinkage and selection operator model. Yang et al. [6,7] studied the forecasting of global horizontal irradiance by exponential smoothing using decompositions, while on another study, they developed the least absolute shrinkage and selection operator model using irradiance very short-term forecasting. Combining a forecasting model with *GSI* is important to get a better result, and forecasting *GSI* by the spatiotemporal pattern recognition method, ANN method, parametric models and decomposition models, has been described in [8–10]. Spatiotemporal pattern recognition and nonlinear principal component analysis (PCA) for global horizontal irradiance forecasting has been proposed as well by Licciardi et al. [8]. Amrouche and Le Pivert [9] have presented an ANN based on daily local forecasting for global solar radiation, describing a novel methodology for local forecasting of daily global horizontal irradiance (GHI). The methodology is a combination of spatial modelling and ANNs algorithm. Wong et al. [10], have presented solar radiation models based on parametric models and decomposition models for predicting the average daily and hourly global radiation, beam radiation and diffuse radiation.

Obtaining *GSI* forecasting for a PV station is also the subject of several studies. In the previous references the *GSI* forecasting was carried out without being influenced by the location of the PV station, but in [11,12] the *GSI* forecasting was very influenced by the Mediterranean location studied. To achieve multi-horizon irradiation forecasting for Mediterranean locations, time series models have been proposed by Paolia et al. [11]. Lorenz et al. [12] presented irradiance forecasting for the power prediction of grid-connected PV systems. In addition to *GSI* forecasting using grid-connected systems there are also other studies that similarly use the grid-connected method but different locations, as explained in [13,14]. Wang et al. [13] studied a short term solar irradiance forecasting model based on an ANN using statistical feature parameters. Another study on *GSI* forecasting using an ANN was performed by Mellit and Pavan [14], which they applied the prediction to a grid-connected PV plant located at Trieste, Italy. The application of a statistical method to detect the motion of cloud structures for surface irradiance is widely used for forecasting *GSI* as in the previous reference, and therefore in [15] optimization and operational validation of the *GSI* forecasting, i.e., short term forecasting of solar radiation using statistical methods to determine cloud motion vector fields have been proposed by Hammer et al. Xiao and Chaovalitwongse [16] have presented optimization models for decomposed nearest neighbor feature selection. Validation of short and medium term operational solar radiation forecasts in the US was studied by Perez et al. [17]. Many of the development models done by other researchers for forecasting *GSI*, tried *GSI* forecasting with statistical methods of stochastic learning and the development of analytical models with time series as in [18,19], i.e., forecasting of global and direct solar irradiance using stochastic learning methods, ground experiments and the national weather service's (NWS) database have been proposed by Marquez et al. [18]. Martin et al. [19] have presented *GSI* forecasting methods based on time series analysis that have been used to predict half daily values of solar irradiance for the next three days. *GSI* forecasting models developed by performing functional and fuzzy approach spatiotemporal model development have been used too, as described in [20,21]. Boata and Gravila [20] have presented a functional fuzzy approach for forecasting daily global solar irradiation and very short term forecasting of the global horizontal irradiance using a spatiotemporal autoregressive model has been proposed by Dambreville et al. [21]. Of the various modeling approaches for forecasting *GSI* in the previous references, namely developing and combining forecasting models for short-term forecasts, [22–24] also describe several methods for forecasting *GSI*, namely using support vector machines and space exponential smoothing models. A review of solar irradiance forecasting methods and a proposition for small-scale insular grids has been provided by Diagne et al. [22]. Short term solar power prediction using a support vector machine has been proposed Zeng et al. [23]. Dong et al. [24] have presented a short

term solar irradiance forecasting method using an exponential smoothing state space model. A comparative empirical study based on a short term wind speed forecasting model has been presented by Ren et al. [25]. Mellit et al. [26] and Farhad et al. [27] have presented ANN models for the prediction of solar radiation and a new bloggers classification approach with a hybrid k -NN and ANN model. Probabilistic solar power forecasting approaches based on k -NN kernel and selection of input parameters to model direct solar irradiance by using ANN have been proposed by Zhang and Wang [28] and López et al. [29]. The aforementioned studies elaborated on GSI forecasting at one target PV station, but have yet to forecast GSI at PV stations surrounded by other PV stations.

This paper presents part of a study in process seeking to estimate and predict one hour or 60 min ahead global solar irradiance at PV stations for energy production. This part focuses on how to predict hourly GSI for the target PV station with the availability of a local database and based on meteorology data. In this study, a new hybrid methodology that combines k -NN modelling and ANN modelling algorithm has been developed. A k -NN-ANN method is used to forecast GSI at the target PV station by means of calculating k -NNs based on the Euclidean distance and then do the testing and training data.

The remainder of the paper is organized as follows: Section 2 describes the modelling and data from a PV station, Section 3 describes the methodology used, i.e., the k -NN and neural network models, while Section 4 presents modelling study cases for very short-term forecasting and their measured errors. Finally, Section 5 presents some concluding remarks.

2. Modelling and Data Description

The GSI measurements were performed continuously every 5 min for four hours. The dataset thus contains four hours of data (from 5:20 a.m. to 8:00 a.m.) collected on 8 June 2012 at nine PV station locations: Station A (0° , 8.3 km), Station B (36° , 10.5 km), Station C (93° , 10.8 km), Station D (140° , 10.5 km), Station E (180° , 10 km), Station F (250° , 4.3 km), Station G (280° , 5.4 km), Station H (310° , 9.1 km) and Station S (0° , 0 km), located in Taipei, Taiwan. For this study, the very short-term forecasts of GSI were only provided by PV Station S which was located in the center. The nearest neighbouring locations surrounding Station S were Station A, Station B, Station C, Station D, Station E, Station F, Station G and Station H (Figure 1).

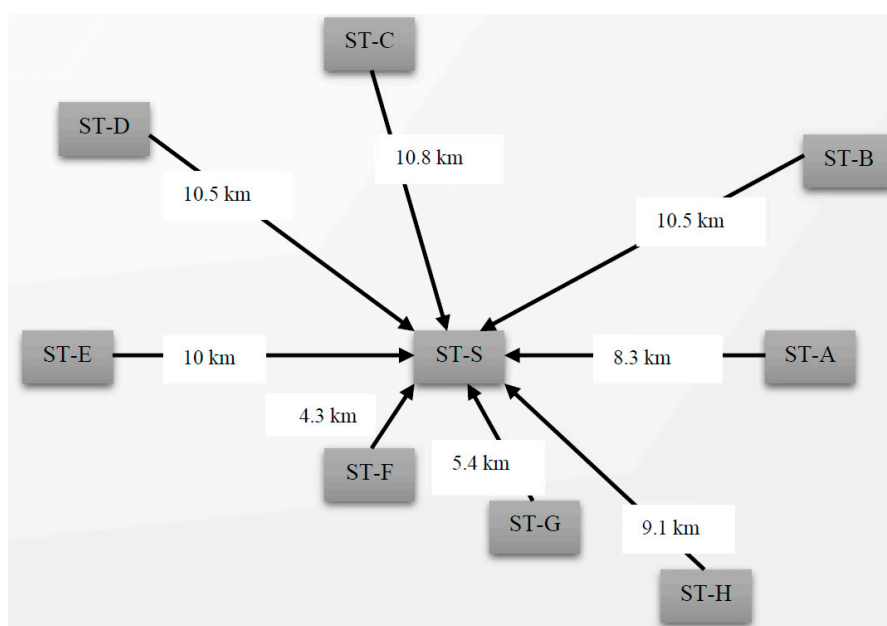


Figure 1. Modelling of Station S which is surrounded by eight other photovoltaic (PV) stations. ST: Station.

Figure 2 illustrates the GSI values as measured at Station S. These monitored data have been used to evaluate the methodology algorithm using k -NN-ANN as the proposed model for GSI forecasting.

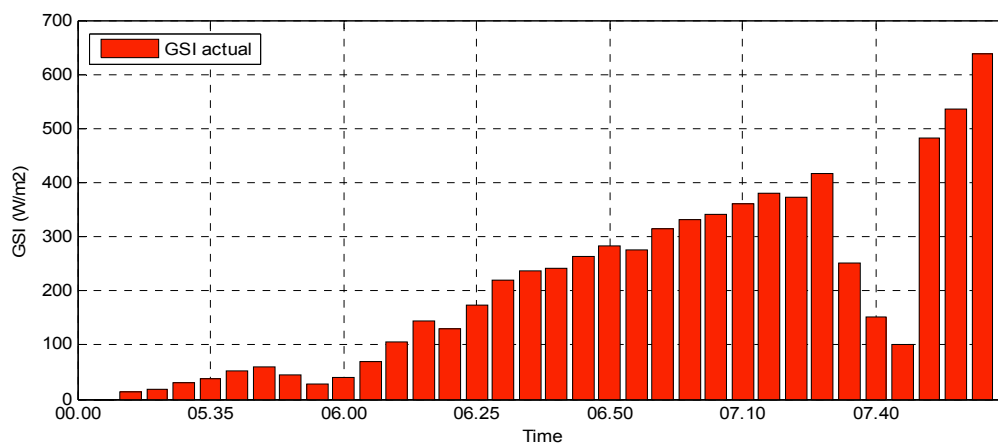


Figure 2. Measured global solar irradiance (*GSI*) at Station S from 5:20 a.m. to 8:00 a.m. on 8 June 2012.

3. *k*-Nearest Neighbor and Artificial Neural Network Model for Forecasting Global Solar Irradiance Based on Meteorological Data

This section explains the basic idea of the construct methodology for *GSI* prediction, namely the *k*-NN-ANN model. In this study, the subject is the central station, which it is surrounded by several other PV stations. The first purpose of this study is the improvement of forecasting results using the *k*-NN method combined with an ANN model method, and the process is then used to predict *GSI* output result of a PV station one hour or 60 min ahead based on meteorological data.

Simulation of the *k*-NN neural networks can be programmed in a few minutes after the recording of the first measurements. For the *GSI* forecasting, *k*-NN-ANN method employs past meteorological data (*GSI*, temperature, humidity, wind speed and direction). The forecast horizon is four hours in 5 min increments. The first step in developing a *k*-NN-ANN method is to develop the database of features that will be used for comparison with the current conditions and the forecast *GSI* a few ahead.

3.1. *k*-Nearest-Neighbors

The *k*-NN method is one of the simplest machine learning algorithm methods. The *k*-NN algorithm is a non-parametric method used for classification and regression. The output depends on whether *k*-NN is used for classification or regression: in used *k*-NN model classification is the value output is a class membership. An object validation is classified by a majority vote of its neighbors, with the object being assigned to the class most common among its *k*-NNs. If *k* = 1, then the object is simply assigned to the class of that single nearest neighbor. In a *k*-NN regression model, the value output is the property value for the object. This value result output training is the average of the values of its *k*-NNs. The *k*-NN model is applied to perform classification of objects based on learning data that were located closest to the object, and the method is considered the simplest among other methods [2]. The main idea is that the *k*-NN algorithm uses a training set for data modelling. Then, the prediction of new points can be the average of the values of its *k*-NNs. The variables employed for modelling very short term forecasting have been described in the previous section. The *k*-NN predict is computed using the features assembled in the matrices in a two-step process. In the first step, we have been calculating the pre-defined distance between the variables in the new dataset (the optimization or the testing sets and training sets) and the features in the previous dataset. For a given set of features $S = \{p^1, \dots, p^n\}$ in the new dataset with lengths N_1, \dots, N_n , the distances of the previous data are calculated. In the second step, choosing *k*-NNs and have *k* smallest distances from training test [26]. The distance *D* is sorted in ascending order, and the first *k* elements $D_S (D_{S,1} \leq D_{S,2} \leq \dots \leq D_{S,k})$ and their associated *k* time stamps $\{\tau_1, \dots, \tau_k\}$ are extracted [2] Equation (1). To find the *k*-NN based on the Euclidean distance, this mathematical equation is used:

$$d(x, y) = \sqrt{\sum_{j=1}^N w_j^2 (x_j - y_j)^2}, j = 1, 2, 3, \dots \quad (1)$$

where d is the number of forecast instances in the optimization set. We can calculate the distance between two scenarios using some distance function $d(x, y)$, where x, y are the matrix scenarios composed of N features $x = \{x_1, \dots, x_N\}, y = \{y_1, \dots, y_N\}$, N is the length of data, and the distance between the current performance and previous condition, w_j is the weight value of the dependent variable members of k -NN (kernel function) and j is the order of the k -NN based on their distance from the current performance condition and which the nearest with used the lowest order ($j = 1, \dots, K$).

k -Nearest Neighbor Modelling

In this research, to obtain the k -NN model forecast using the algorithm model proposed described above, several parameters need to be specified and the k -NN modelling process can be divided into five calculation stages. The procedure of k -NN for regression is as follows:

- (1) The matrix scenarios composed modelling stage which includes form d -dimensional feature vectors C or nxy from the historical data $x: c = [c_1, c_2, \dots, c_p]$ and $nxy = [nxy_1, nxy_2, \dots, nxy_p]$ or $x: C = nxy = [x_t, x_{t-1}, \dots, x_{t-d+1}]$; Their corresponding successors are denoted as x^h . They are given two pieces point c and nxy in a space vector of n -dimensional $c (c_1, c_2, \dots, c_n)$ and $nxy (nxy_1, nxy_2, \dots, nxy_n)$.
- (2) The distance calculate vector stage which includes form n -dimensional distance vector $ds_j^i(c, nxy) = D_i$ for each testing vector C or nxy by calculating the Euclidean distance between $ds_j^i(c, nxy) = D_i$ and the remaining:

$$C : ds_j^i(c, nxy) = \left\{ \|D_i - D_j\| \right\}, j \neq i \quad (2)$$

where D_i is the value of the dependent variable historical data set and D_j is the value of the dependent the nearest neighbor based on distance:

$$ds_j^i(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^i - nxy_{pj}^i)^2}, p = 1, 2, 3, \dots, n \quad (3)$$

where $ds_j^i(c, nxy)$ is value of the dependent variable GSI , nxy_{pj}^i is the position coordinat PV station (magnitude of nearest neighbors) c_p^i is the d -dimensional feature vectors, The index j is the current condition, i is the historical dataset, $k = K$ the number of elements in the nearest neighbors, q is the historical dataset and where x, y are scenarios composed of k features and p is the number of the k -NN based on their distance from the current condition (j) in which the nearest have the lowest order ($p = 1, \dots, k$).

- (3) Select the value distance of the k -NN stage, which includes the sort $ds_j^i(c, nxy)$ in ascending order and select the first K entries as the nearest neighbors $D_k, k \in \{1, \dots, K\}$.
- (4) Select the best value of k used in modelling the k -NN stage, because a high k value will reduce the effect of noise on the classification, but it will make the boundaries between each classification becomes increasingly blurred. Form a kernel function:

$$K(j) = \frac{1/j}{\sum_{j=1}^K 1/j} \text{ or } k_i = 1/\text{dis}(i) \quad i = 1, 2, 3, 4, 5 \quad (4)$$

where $K(j) = k_i$ is the value k -NN (kernel function), i is the historical data set and the index j is the number of the k -NN based on their distance from the current condition (i) in which the nearest have the lowest order ($j = 1, \dots, K$), and K is the length data sets, and the distance between the current and previous condition [25].

- (5) Calculate the final estimation stage, using Equation (5):

$$\text{sumd} = \sum_{j=1}^k K_j x_j^h \quad (5)$$

where x_j^h is the magnitude of nearest neighbor j , j is the order of the nearest neighbors based on their distance from the current condition (h) in which the nearest have the lowest order ($j = 1, \dots, k$), sumd is the kernel function, k is the length of data sets, and K_j is the kernel function.

3.2. Artificial Neural Network

ANN is a mathematical method inspired by the structure and information processing of biological neural networks. ANNs are intelligent systems that have the capacity to learn, memorize and create relationships among data [30]. ANN model is a combination of pattern recognition, deductive reasoning and numerical computations to simulate learning in the human brain. ANN consists of an interconnected many groups namely neurons, and its main task is processes information using a connection approach to computation. ANN consists of an interconnected many groups namely neurons, and its main task is processes information using a connection approach to computation. ANN models have been used to predict solar radiation data [24]. The methodology is a promising alternative to traditional approaches for forecasting *GSI*, especially in cases where radiation measurements are not readily available. ANN models fundamentally comprise multiple connected neurons and nodes. The neural networks are considered as the member of the non-parametric techniques which are usually used for estimation and classification [25]. The neurons have five basic components, i.e., input, weight-bias, threshold, summing junction and output, as illustrated in Figure 3. Neurons are arranged in three layers which consist of input, hidden and output.

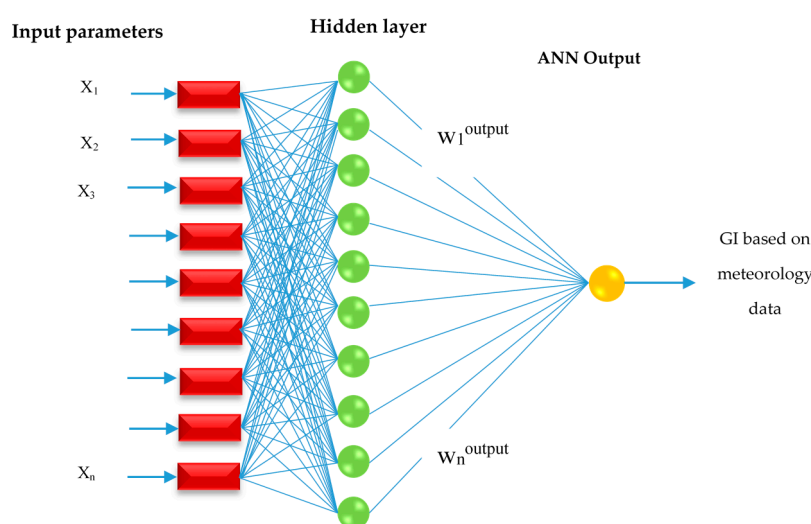


Figure 3. Basic structure architecture of a simple artificial neuron model [27].

Artificial Neural Network Modelling

The ANN algorithm model can be divided into four step: (1) the design model and input pattern step which includes the choice of the ANN model, the number of its layers ANN, the number of neuron groups in each layer ANN, its inputs and outputs, the choice of training set and validation set samples which use k -NN design; (2) the training set and testing set ANN forecasting based on meteorology data are presented to the ANN and the weights layer are adjusted accordingly till a predetermined condition is more better; (3) the simulation test passed step, in which the result ANN forecasting model is tested using measurement data at nine station PV which are not treated during the training step; and finally (4) the performance evaluation stage.

Layered perceptron ANN model with different topology designs were considered in order to obtain the best mapping between the ANN algorithm inputs and outputs. The proposed model in this research is used to predict the *GSI* values for the next hours or one hour ahead, forecasting based on meteorology data and the *GSI* data from other PV stations. The *GSI* forecasting data are

normalized to the limit of [0, 1] to avoid neuron groups saturation during the learning progress. The number of neuron groups in the first hidden layer is 10 and we have to consider an ANN with 1475 inputs for *GSI* forecasting. The neuron groups active function of hidden and output layer are tansig and purelin, respectively. Performance test data set is used to evaluate *GSI* forecast accuracy of the ANN algorithm. The ANN models inputs acts directly on the training period duration for *GSI* forecasting [9]. To choose the best design modeling for forecasting, we investigated the role of meteorology data in the ANN model effectivity and then we researched the impact of the data mode on the training good performance and on the forecasting accuracy.

In the first reason, the ANN models will learn about the global solar irradiance profile of the target location. At the end of the learning process, the *k*-NN-ANN model will be able to give the *GSI* forecast results at the target PV station based on meteorology data. When compared to actual data, ANN also presents the whole *GSI* values for the four-hour window, giving ANN the possibility to learn the existing relationship in between them and to develop an idea about *GSI* evolution during the time window. In the second reason, the entire *GSI* data are used for very short term forecasting using the *k*-NN-ANN model based on meteorological data at the target PV station, which is surrounded by eight other PV stations. The proposed *k*-NN-ANN model has to forecast *GSI* values for the next hour or 60 min ahead by taking into account the forecasting data based on meteorology data.

The specific feature of this *k*-NN-ANN model is that measured data and hourly weather meteorology data (temperature, humidity, global irradiance wind speed and direction) forecasts for the nearest eight surrounding stations are used as input data information, as explain by Figure 4. To model the available relationship between *k*-NN model prediction and the actual values of *GSI* at the pivot location, an ANN based model is used. The training of the ANN models can be programmed in one hours ahead after the recording of first measurements.

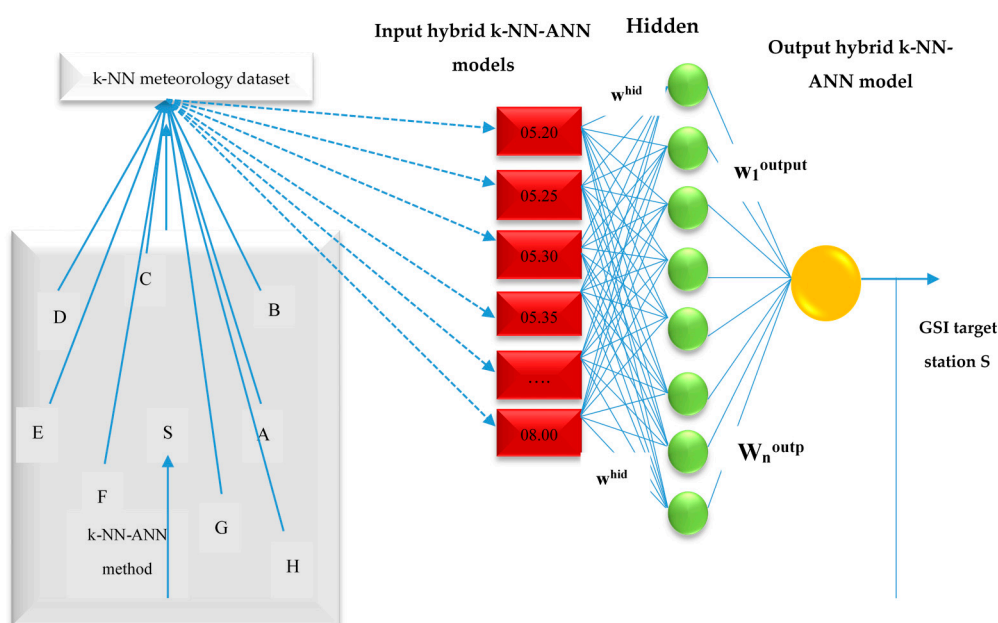


Figure 4. The proposed forecasting model using a *k*-nearest neighbor and artificial neural network (*k*-NN-ANN) model.

3.3. *k*-Nearest Neighbor and Artificial Neural Network Modelling

The forecasting problem in this research was forecasting based on meteorological data study in progress to forecast *GSI* at a PV station for energy production by one hour or 60 min ahead by using the hybrid *k*-NN-ANN model. The procedure of *k*-NN-ANN model for *GSI* forecasting based on meteorological data has the following steps:

Step 1: Calculate the distance of each data parameter:

$$ds_j^{GI}(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^{GI} - nxy_{pj}^{GI})^2}, \quad p = 1, 2, \dots, n \quad j = 1, 2, \dots, n \quad (6)$$

$$ds_j^{Ta}(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^{Ta} - nxy_{pj}^{Ta})^2}, \quad p = 1, 2, \dots, n \quad j = 1, 2, \dots, n \quad (7)$$

$$ds_j^{Ho}(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^{Ho} - nxy_{pj}^{Ho})^2}, \quad p = 1, 2, \dots, n \quad j = 1, 2, \dots, n \quad (8)$$

$$ds_j^{Ws}(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^{Ws} - nxy_{pj}^{Ws})^2}, \quad p = 1, 2, \dots, n \quad j = 1, 2, \dots, n \quad (9)$$

$$ds_j^{Wd}(c, nxy) = \sqrt{\sum_{q=1}^k (c_p^{Wd} - nxy_{pj}^{Wd})^2}, \quad p = 1, 2, \dots, n \quad j = 1, 2, \dots, n \quad (10)$$

and the weight distance is:

$$GI_x = \frac{1}{\sum d_1^2} \cdot GI_i + \frac{1}{\sum d_2^2} \cdot GI_i + \frac{1}{\sum d_3^2} \cdot GI_i + \frac{1}{\sum d_4^2} \cdot GI_i + \frac{1}{\sum d_5^2} \cdot GI_i \quad (11)$$

$$Ta_x = \frac{1}{\sum d_1^2} \cdot Ta_i + \frac{1}{\sum d_2^2} \cdot Ta_i + \frac{1}{\sum d_3^2} \cdot Ta_i + \frac{1}{\sum d_4^2} \cdot Ta_i + \frac{1}{\sum d_5^2} \cdot Ta_i \quad (12)$$

$$Ho_x = \frac{1}{\sum d_1^2} \cdot Ho_i + \frac{1}{\sum d_2^2} \cdot Ho_i + \frac{1}{\sum d_3^2} \cdot Ho_i + \frac{1}{\sum d_4^2} \cdot Ho_i + \frac{1}{\sum d_5^2} \cdot Ho_i \quad (13)$$

$$Ws_x = \frac{1}{\sum d_1^2} \cdot Ws_i + \frac{1}{\sum d_2^2} \cdot Ws_i + \frac{1}{\sum d_3^2} \cdot Ws_i + \frac{1}{\sum d_4^2} \cdot Ws_i + \frac{1}{\sum d_5^2} \cdot Ws_i \quad (14)$$

$$Wd_x = \frac{1}{\sum d_1^2} \cdot Wd_i + \frac{1}{\sum d_2^2} \cdot Wd_i + \frac{1}{\sum d_3^2} \cdot Wd_i + \frac{1}{\sum d_4^2} \cdot Wd_i + \frac{1}{\sum d_5^2} \cdot Wd_i \quad (15)$$

where $ds_j^{GI}(c, nxy)$ is the Euclidan distance GSI

$ds_j^{Ta}(c, nxy)$ is the Euclidan distance temperature

$ds_j^{Ho}(c, nxy)$ is the Euclidan distance humidity

$ds_j^{Ws}(c, nxy)$ is the Euclidan distance wind speed

$ds_j^{Wd}(c, nxy)$ is the Euclidan distance wind direct

GI_x is the weight distance at GSI

Ta_x is the weight distance at temperature

Ho_x is the weight distance at humidity

Ws_x is the weight distance at wind speed

Wd_x is the weight distance at wind direct

GI is the global irradiance

Ta is the temperature

Ho is the humidity

Ws is the wind speed

Wd is the wind direct

C_p^{GI} is the d dimensional feature vectors GSI

C_p^{Ta} is the d dimensional feature vectors temperature

C_p^{Ho} is the d dimensional feature vectors humidity

C_p^{Ws} is the d dimensional feature vectors wind speed

C_p^{Wd} is the d dimensional feature vectors wind direct

nxy_{pj}^i is the position coordinat PV station (magnitude of nearest neighbors)

C_p^i is the d -dimensional feature vectors

sumd is the kernel function:

i is the historical data set

j is the current condition, i is the historical dataset, $k = K$ the number of elements in the nearest neighbors

q is the historical dataset and where x, y are scenarios composed of k features and p is the number of the k -NN based on their distance from the current condition (j) in which the nearest have the lowest order ($p = 1, \dots, k$).

Step 2: Calculate number of nearest neighbors:

$$k_{i+n} = \frac{1}{\text{dis}(i+n)} \quad i = 1, 2, 3, \dots, n = 130 \quad (16)$$

where k is the distance nearest neighbor, i is the historical dataset and n is the total number of features ($i = 1, 2, \dots, n$)

Step 3: Calculate the final estimation as:

$$\text{sumd} = \sum_{j=1}^k K(j)x_j^h \quad (17)$$

where x_j^h is the value magnitude of the k -NN j , j is the order of the nearest neighbors based on distance (h) in which the nearest have the lowest order ($j = 1, \dots, k$), sumd is the kernel function, k is the length of data sets, and the distance between the current and previous condition and K_j is the kernel function.

Step 4: Training data and testing data using ANN method are obtained from the following steps:

- The final estimation of GSI from k -NN method is into training sets, validation sets, and test sets for ANN model.
- Architecture model and training sets parameters (create and configure the neural network, initialize the weights layer and biases layer) are selected
- Model GSI forecast using the training set are run
- Model forecasting is validated using the validation set
- Step b–e are repeated using different architectures model forecasting and training set parameters
- The optimum model is chosen and inserted into training process using data from the training set and validation set
- The ultimate forecasting model is assessed using the test set

Step 5: Perform predictions for future GSI data at the Station S with k -NN-ANN

The forecast model designed for k -NN is described in Figure 5a. If the validation test for forecasting the GSI is successful and get more better the result, the forecast model could perform its designed function, otherwise one or more changes should be made during the previous process. The procedures are shown in Figure 5b. GSI forecasting process using hybrid k -NN-ANN algorithm can also be divided into two process: (1) k -NN modelling to determine the d -dimensional feature and n -

dimensional distance for input data at the ANN process; (2) the ANN modelling for GSI forecasting. The procedures are displayed in Figure 5c. In this research, *k*-NN and ANN structure construction was programmed using MATLAB (R2013a) programming. MATLAB is been provided with some ANN tools.

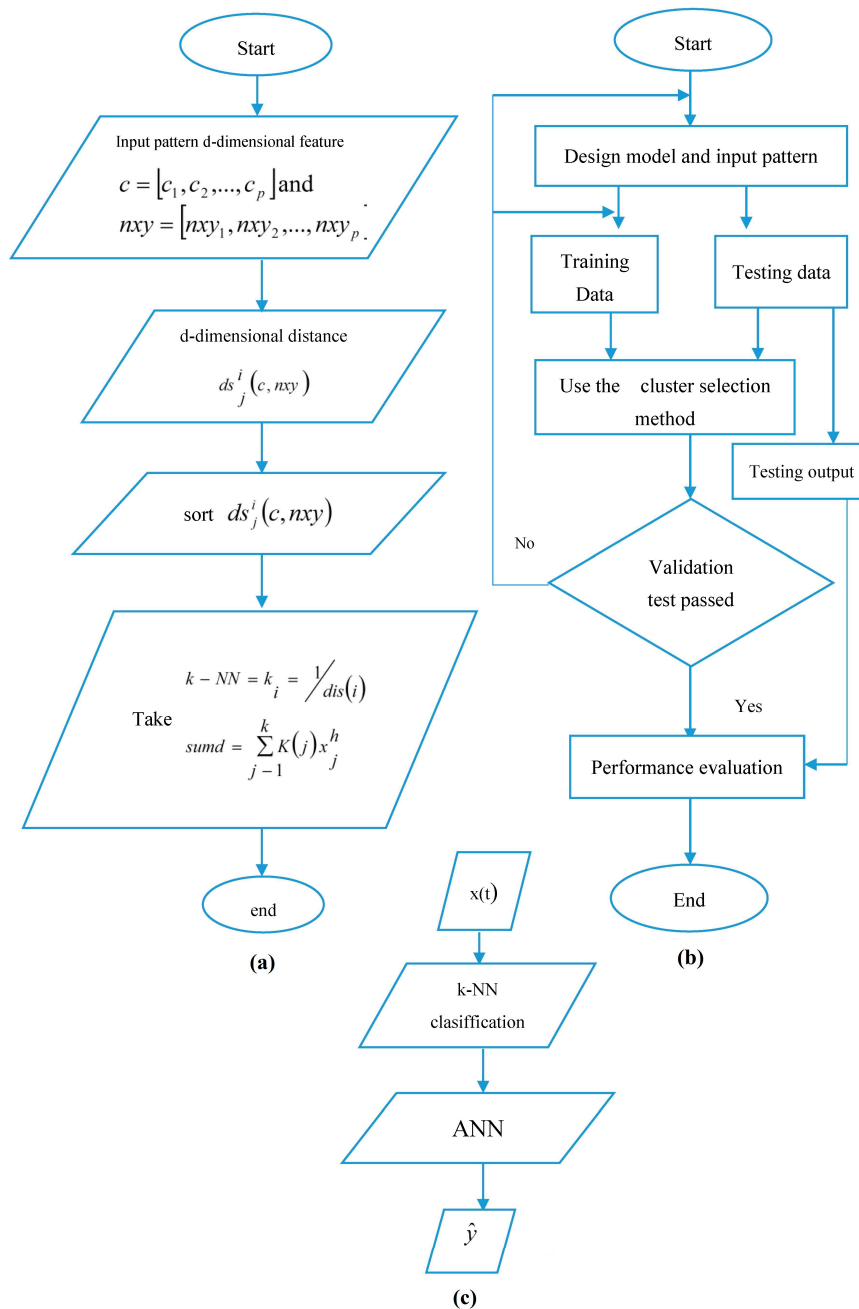


Figure 5. Flowchart *k*-NN-ANN method very short term forecasting: (a) *k*-NN process; (b) ANN process; and (c) hybrid *k*-NN-ANN model.

3.4. Normalization

The normalization of data input is very important to obtain good results in the ANN method [9]. In this research for analysis, it needs normalization data for training process the GSI forecasting 60 min ahead, as defined [31], which can be calculated by Equation (18):

$$GSI_{dnorm}^n = \frac{GSI_d^n - GSI_{min}}{GSI_{max} - GSI_{min}} \tag{18}$$

Let us denote GSI_{dnorm}^n and GSI_d^n be normalized target feature at frame index n and the d - GSI output, respectively. Let us also denote GSI_{max} and GSI_{min} be the maximum value and minimum values of the GSI , respectively.

4. Results and Discussion

This section discusses the result of k -NN-ANN model used to forecast future GSI by using a one hour ahead forecasting procedure. The procedure is described in Figure 5. We tested our model using previously described databases based on meteorology data, i.e., wind direction, wind speed, GI , temperature, and humidity, during the 1 h or 60 min ahead process.

Data

Using the k -NN-ANN model, it is expected that a valid GSI forecast result will be produced. The design forecast is divided into two stages:

- The first stage is calculating d -dimensional feature and n -dimensional distance based on the Euclidean (k -NN model) every hour for all of the PV stations. Data shown in Table 1 are the order parameters of each PV station, i.e., angle, distance and position coordinates. The table summarizes all possible combinations of variables to be considered as an input for k -NN method.
- The second stage uses the ANN based on the assumption that the existing data input is a combination of the results obtained from the k -NN model. The research model proposed in this study seeks to estimate and predict a PV station production 60 min ahead, which position is located at the center and surrounded by eight other PV stations.

Table 1. Data position and coordinates of the PV stations.

No.	Station	Angle (°)	Distance (d)	Coordinate (x_i, y_i)
1	A	0	8.3	(8.3, 0)
2	B	36	10.5	(8.5, 6.2)
3	C	93	10.8	(−0.6, 10)
4	D	140	10.5	(−8, 6.7)
5	E	180	10	(−10, 0.1)
6	F	250	4.3	(−5.5, −5)
7	G	280	5.4	(1, −5.3)
8	H	310	9.1	(5.8, −7)
9	S	0	0	(0.1, 0.1)

The d -dimensional feature and n -dimensional distances based on the Euclidean (k -NN model) every hour for all of PV station are shown in Figure 6. From the simulation results using the k -NN method based on meteorological data consisting of global irradiance, temperature, humidity, wind speed and wind direction values, respectively. All of them are used as pre-processing data in the ANN method, which can be calculated by the polynomial Equation (19):

$$f(x, y) = p_{00} + p_{10}x + p_{01}y + p_{20}x^2 + p_{11}xy + p_{02}y^2 \quad (19)$$

where $f(x, y)$ is the GSI value of the k -NN method and variable x, y is the coordinate value's position of the PV station.

Example Figure 6a can be produced by polynomial Equation (20):

$$f(x, y) = 10.74 + 0.05105x + 0.01365y + (-0.0014)x^2 + 0.0015xy + 0.002923y^2 \quad (20)$$

Example Figure 6b can be produced by polynomial Equation (21):

$$f(x, y) = 520.4 - 1.678x - 0.7583y + (-0.3799)x^2 + 0.2108xy + 0.07159y^2 \quad (21)$$

In which the polynomial equation above is the result of the simulation on 5:20 a.m. for Figure 6a and 8:00 a.m. hours for Figure 6b using the *k*-NN method for Station S. Moreover that polynomial equation also can be implemented for the other PV stations.

To validate the proposed method, *GSI* data of 60 min ahead of a PV station has been calculated using the process described in Section 3. The results are then compared with the actual data of the *GSI* at target PV station as described in Section 4. Table 2 shows the optimal *k*-NN parameters for the *GSI* forecast with meteorology data every hours from 5:20 a.m. to 8:00 a.m. on 8 June 2012.

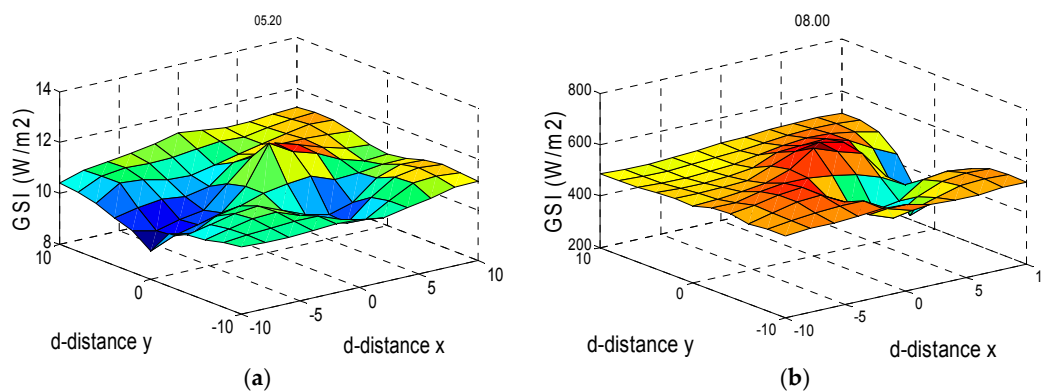


Figure 6. (a) *d*-Distance PV station based on the Euclidean (*k*-NN model) in time 5:20 a.m. on 8 June 2012; and (b) *d*-Distance PV station based on the Euclidean (*k*-NN model) in time 8:00 a.m. on 8 June 2012.

Table 2. Optimal *k*-NN parameters for the *GSI* forecast with meteorology data.

Time	Station (<i>x, y</i>)*								
	ST-S (0.1, 0.1)	ST-A (8.3, 0)	ST-B (8.5, 6.2)	ST-C (-0.6, 10)	ST-D (-8, 6.7)	ST-E (-10, 0.1)	ST-F (-5.5, -5)	ST-G (1, -5.3)	ST-H (5.8, -7)
5:20	11	11	11	11	10	10	10	11	11
5:25	15	15	15	15	15	15	15	15	15
5:30	22	22	22	22	22	22	22	22	22
5:35	27	29	27	25	24	25	27	28	30
5:40	37	39	37	33	33	35	38	39	41
5:45	42	44	44	42	41	40	41	42	43
5:50	33	35	38	38	35	31	30	31	31
5:55	32	32	31	31	31	32	32	32	32
6:00	40	41	36	33	35	40	44	44	46
6:05	59	61	53	45	48	57	65	67	70
6:10	77	82	76	67	66	72	79	83	88
6:15	101	106	100	91	90	95	103	107	111
6:20	139	132	113	110	127	148	160	155	155
6:25	169	152	133	140	165	189	195	183	178
6:30	179	158	153	172	194	204	197	181	170
6:35	197	186	184	196	207	211	206	197	191
6:40	216	200	202	221	234	235	225	213	203
6:45	230	217	219	236	246	246	237	226	218
6:50	240	229	233	248	256	254	245	235	228
6:55	252	238	243	263	272	269	256	245	235
7:00	269	256	260	278	287	285	274	263	254
7:05	292	275	275	294	310	314	304	290	280
7:10	308	290	290	310	326	330	320	305	295
7:15	324	302	302	327	347	351	338	321	307
7:20	345	322	320	345	366	373	361	343	330
7:25	357	334	334	360	381	385	372	354	340
7:35	370	351	354	379	394	395	380	365	351
7:40	362	331	334	371	397	400	380	355	336
7:45	360	323	327	372	402	405	380	351	327
7:50	344	316	328	367	385	378	352	329	310
7:55	427	414	414	429	440	443	436	426	418
8:00	475	457	450	467	486	497	492	478	469

Note: * (*x, y*) are the coordinate position values of the PV stations.

The k -NN-ANN modelling approach is unique since the neural network algorithm is firstly developed using the k -NN model based on meteorology data. Results from the k -NN model are then divided it into two sets of data: the training and the validation set. After the simulation test, the ANN model based on the k -NN results will be ready to use in 60 min ahead GSI forecasting. As explained previously, the optimization parameters for GSI forecasting were obtained from a dataset based on meteorology data. The calculated values of GSI were then compared with measured values (GSI) of each station: Station S, Station A, Station B, Station C, Station D, Station E, Station F, Station G, Station H. Mean absolute bias error ($MABE$) and root-mean-square error ($RMSE$) were used as error statistical indicators. The estimation from the k -NN-ANN model was then compared with the k -NN model and mean average of meteorological data prediction.

For the comparison, k -NN-ANN model performed only 500 iterations for each learning period. Figure 7 shows the all the data for the three subsets: (a) training dataset; (b) validation dataset; and (c) test dataset. In this figure the plots show for GSI versus the time. For the ANN algorithm is initially constructed for the training based on the k -NN model data, and after this, its training is periodically as the database expands over time.

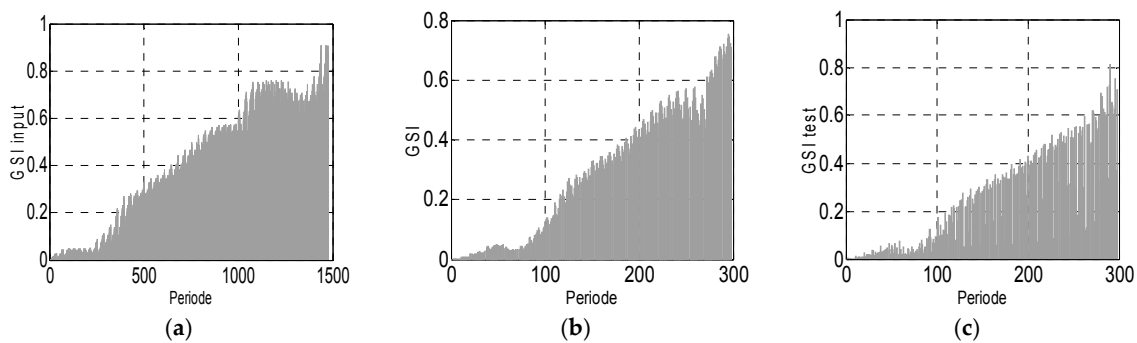


Figure 7. (a) GSI for the training set; (b) GSI for the validation set; and (c) GSI for the testing set.

The k -NN model database provides pretraining data, and then the database divided it into two sets of data; the training data set and the validations data set, after the validation test, the k -NN-ANN based model will be ready to use for forecasting the GSI , and the the accuracy of the model can be determined using the testing set data that has been gained from the training data process. For the present forecast application program, the ANN model has to start learning based on the training data set, and subsequently construct and sharpen the knowledge while ensuring the continuity of its task. For the process of testing the accuracy of the forecasting model that has been gained from the training process using backpropagation method. The amount of testing data used was 90% of the total. Note that for this validation, the k -NN-ANNs were trained only while data processing is done on the training data patterns with a target error of 0.001, learning rate of 0.1 and a maximum of 50 epochs.

Figure 8 shows a comparison of the two methods for the GSI data between actual data and k -NN-ANN method, where: (a) there is a better between actual and forecasting data for very short term GSI forecasting at Station target S; and (b) shows the normalized GSI curve using the k -NN-ANN method.

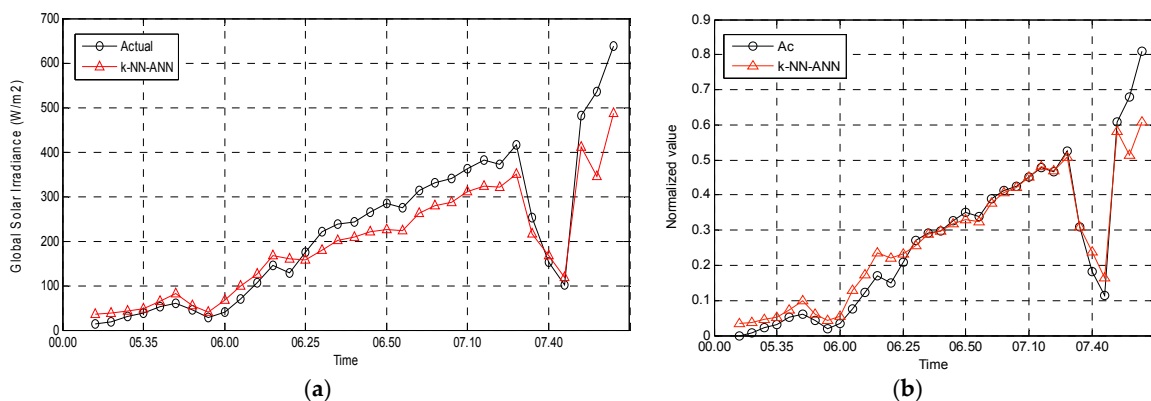


Figure 8. Comparison of the two methods for the GSI data: (a) very short term GSI forecasting using the *k*-NN-ANN model versus actual data for station S; and (b) the normalized GSI curve using the *k*-NN-ANN method.

Figure 9 illustrates the comparison of GSI forecasting in a four hour window (5:20 a.m.–8:00 a.m.) on 8 June 2012, based on the *k*-NN-ANN model, actual data, and the *k*-NN method. The result shows that very short term forecasting simulation using *k*-NN-ANN method during a four hour window gives better results compared to the *k*-NN method.

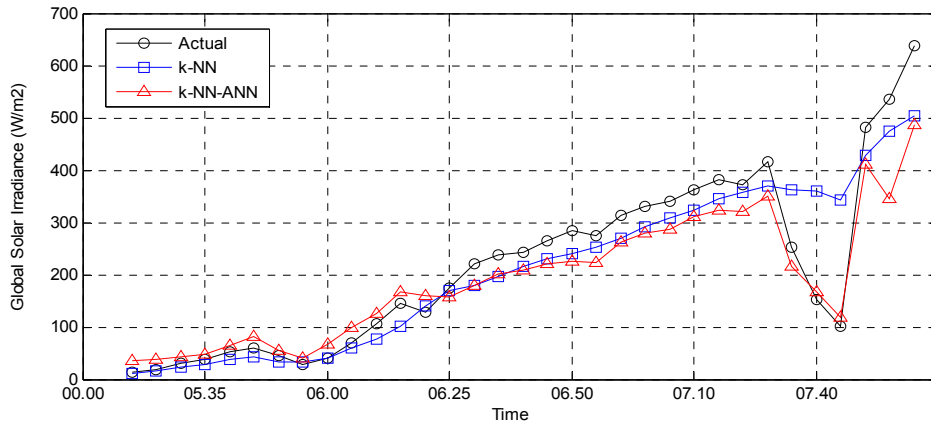


Figure 9. Very short term GSI forecasting in 5-h window based *k*-NN-ANN model, actual data, and *k*-NN method.

Figure 10 illustrates a very short-term (60 min ahead) GSI forecast using *k*-NN-ANN and its comparison with actual data. It is evident that the *k*-NN-ANN model is in a good agreement with the measured data at the object station.

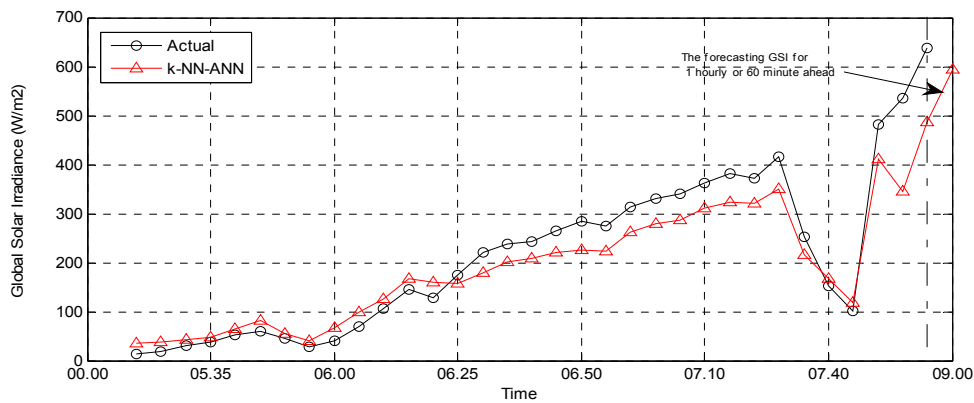


Figure 10. Very short term (60 min ahead) GSI forecast using *k*-NN-ANN versus actual data.

To evaluate the performance of the models, a statistical error measurement was used in the experiment, namely the *MABE*, and *RMSE*. To evaluate the accuracy of each method to forecast the GSI values, *MABE*, and *RMSE* coefficient between results of *k*-NN-ANN and actual ground measurements were calculated. These statistical error indicators validation forecasting are calculated according to Equations (22) and (23) and the results statistical error are shown in Table 3.

Table 3. Error statistical indicators of the GSI forecasting models. *MABE*: mean absolute bias error; *RMSE*: root-mean-square error.

Model Error Indicators	<i>k</i> -NN	<i>k</i> -NN-ANN
<i>MABE</i> (W/m ²)	44	42
<i>RMSE</i> (W/m ²)	251	242

MABE is calculated according to Equation (21):

$$MABE = \frac{1}{N} \sum_{i=1}^N \left(\left| G_{f,i} - G_{m,i} \right| \right) \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (22)$$

RMSE is calculated according to Equation (22):

$$RMSE(k) = \left[\frac{1}{N} \sum_{i=1}^N e^2(t+k|t) \right]^{1/2} \quad (23)$$

where $e(t)$ is the forecasting data and $k(t)$ is the measured (observed) data:

$$p(t+k|t) = P(t+k) - P(t+k|t)$$

where $G_{f,i}$ is forecasted value *GSI* and $G_{m,i}$ is measured value *GSI*, ($i = 1, 2, \dots, N$), N is the number of the *GSI* data, i is the number index variations.

Figure 11 illustrates the *MABE* and *RMSE* coefficients for the actual data and the very short term (60 min ahead) forecasting performance using k -NN and the k -NN-ANN model. The error statistical indicators of the k -NN model are *MABE* 44 W/m² and *RMSE* 251 W/m². On the other hand, the error statistical indicators for the proposed model (k -NN-ANN model) are *MABE* 42 W/m² and *RMSE* 242 W/m². Note that the highest *RMSE* was 191 (W/m²) during the thirty-two period, while the lowest a value was found to be 9 (W/m²) for the same location. It is evident that k -NN-ANN model displays better predictions than the k -NN model.

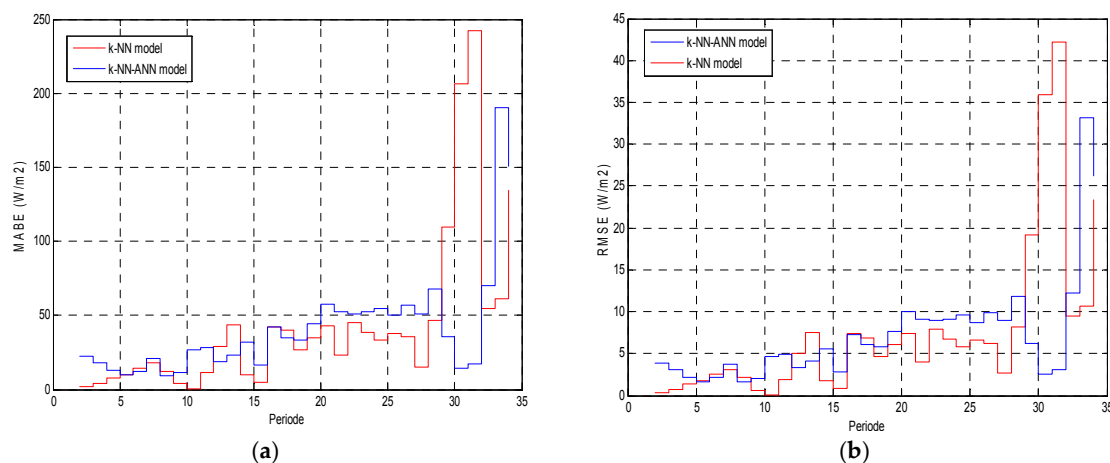


Figure 11. (a) *MABE* coefficients between actual data and *GSI* forecasts using the k -NN-ANN model on the test data set; and (b) *RMSE* coefficients between actual data and *GSI* forecasts using the k -NN-ANN model on the test data set.

5. Conclusions

A new methodology for very short term (60 min ahead) *GSI* forecasting of a target PV station has been introduced. In this work we propose a novel methodology for *GSI* forecasting using a combination of k -NN modelling and an ANN. The model estimates and predicts the *GSI* profiles of PV stations in very short term (60 min ahead) based on hourly meteorology data from eight surrounding PV stations. The following conclusions can be drawn from this research:

- A different formulation for very short term *GSI* forecasting using k -NN-ANN modelling based on meteorology data is proposed. The proposed model attempting to shape the patterns of a polynomial equation shows that the proposed model forecasting is more better. The variable meteorology data weather is of great very importance and affects the resulting *GSI* forecasting output.
- The new model proposed in this study is a combination of k -NN modelling and an ANN model. The model is employed to forecast *GSI* data for very short term period (60 min ahead) based on

meteorology data. This research concerns how to predict *GSI* data at a target PV station, which is surrounded by eight other PV stations. The study also considers the availability of a local measured database. It clearly shows that the *GSI* forecasting using a different *k*-NN-ANN model for every hour based on meteorological data giving a better result output, which means the *GSI* forecasting largely depends on variable meteorological data, where the meteorology data variables consist of *GI*, wind speed, wind direct, humidity and temperatures.

- This paper utilises *k*-NN-ANN modelling to determine the *d*-dimensional features and *n*-dimensional distances. The results demonstrate that the results of the *k*-NN-ANN model are closely matched with the actual data, and are better than data obtained from *k*-NN models.

The novelty of this article is to predict *GSI* at a PV station which position is at the center surrounded by eight other adjacent PV stations. The proposed model is able to learn the characteristics of meteorology weather data for the past four hours and use the data as model input. In this paper, the proposed *k*-NN-ANN model has better approximation compared to *k*-NN model. The very short term forecast evaluations of the *GSI* using the *k*-NN-ANN model are performed for only four hours and the results show that the *k*-NN-ANN method is better than the *k*-NN method. The error statistical indicators of the *k*-NN model are 44 W/m² for the *MABE* and 251 W/m² for the *RMSE*. On the other hand, the error statistical indicators for the proposed model (*k*-NN-ANN model) are 42 W/m² (*MABE*) and 242 W/m² (*RMSE*). We noted that the highest *RMSE* was 191 (W/m²) during the thirty-two hour period, while the lowest value was found to be 9 (W/m²) for the same location. The performance of the proposed *k*-NN-ANN method is more better compared with the *k*-NN model. The proposed *k*-NN-ANN model can therefore be used effectively to forecast very short term *GSI* data while giving closer result outputs and better matches with actual measured data.

Author Contributions: Chao-Rong Chen and Unit Three Kartini conceived and designed the experiments; Chao-Rong Chen performed the experiments; Chao-Rong Chen and Unit Three Kartini analyzed the data; Chao-Rong Chen and Unit Three Kartini contributed reagents/materials/analysis tools; Unit Three Kartini wrote the paper.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Yang, H.T.; Yang, P.C.; Huang, C.L. Optimization of Unit Commitment Using Parallel Structures of Genetic Algorithm. In Proceedings of the 1995 International Conference on Energy Management and Power Delivery, Singapore, 21–23 November 1995.
2. Chiang, C.-T.; Lee, Y.-S.; Li, X.R.; Liao, C.-C. A RSCMAC Based Forecasting for Solar Irradiance from Local Weather Information. In Proceedings of the WCCI 2012 IEEE World Congress on Computational Intelligence, Brisbane, Australia, 10–15 June 2012.
3. Hocaoglu, F.O. Stochastic approach for daily solar radiation modeling. *Sol. Energy* **2011**, *85*, 278–287.
4. Pedro, H.T.C.; Coimbra, C.F.M. Nearest-neighbor methodology for prediction of intra-hour global horizontal and direct normal irradiances. *Renew. Energy* **2015**, *80*, 770–782.
5. Salcedo-Sanz, S.; Casanova-Mateo, C.; Munoz-Mari, J.; Camps-Valls, G. Prediction of Daily Global Solar Irradiation Using Temporal Gaussian Processes. *IEEE Geosci. Remote Sens. Lett.* **2014**, *11*, 1936–1940.
6. Yang, D.; Sharma, V.; Ye, Z.; Lim, L.I.; Zhao, L.; Aryaputera, A.W. Forecasting of global horizontal irradiance by exponential smoothing using decompositions. *Energy* **2015**, *81*, 111–119.
7. Yang, D.; Ye, Z.; Lim, L.H.I.; Dong, Z. Very short term irradiance forecasting using the lasso. *Sol. Energy* **2015**, *114*, 314–326.
8. Licciardi, G.A.; Dambreville, R.; Chanussot, J.; Dubost, S. Spatiotemporal Pattern Recognition and Nonlinear PCA for Global Horizontal Irradiance Forecasting. *IEEE Geosci. Remote Sens. Lett.* **2015**, *12*, 284–288.
9. Amrouche, B.; Le Pivert, X. Artificial neural network based daily local forecasting for global solar radiation. *Appl. Energy* **2014**, *130*, 333–341.
10. Wong, L.T.; Chow, W.K. Solar radiation model. *Appl. Energy* **2001**, *69*, 191–224.
11. Voyant, C.; Paolia, C.; Musellia, M.; Niveta, M.-L. Multi-horizon Irradiation Forecasting for Mediterranean Locations Using Time Series Models. *Energy Procedia* **2014**, *57*, 1354–1363.
12. Lorenz, E.; Hurka, J.; Heinemann, D.; Beyer, H.G. Irradiance Forecasting for the Power Prediction of Grid-Connected Photovoltaic Systems. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* **2009**, *2*, 2–10.

13. Wang, F.; Mi, Z.; Su, S.; Zhao, H. Short-Term Solar Irradiance Forecasting Model Based on Artificial Neural Network Using Statistical Feature Parameters. *Energies* **2012**, *5*, 1355–1370.
14. Mellit, A.; Pavan, A.M. A 24-h forecast of solar irradiance using artificial neural network: Application for performance prediction of a grid-connected PV plant at Trieste Italy. *Sol. Energy* **2010**, *84*, 807–821.
15. Hammer, A.; Heinemann, D.; Lorenz, E.; Luckehe, B. Short-Term Forecasting of Solar Radiation: A Statistical Approach Using Satellite Data. *Sol. Energy* **1999**, *67*, 139–150.
16. Xiao, C.; Chaovalitwongse, W.A. Optimization Models for Feature Selection of Decomposed Nearest Neighbor. *IEEE Trans. Syst. Man Cybern. Syst.* **2016**, *46*, 2168–2216.
17. Perez, R.; Kivalov, S.; Schlemmer, J.; Hemker, K., Jr.; Renné, D.; Hoff, T.E. Validation of short and medium term operational solar radiation forecasts in the US. *Sol. Energy* **2010**, *84*, 2161–2172.
18. Marquez, R.; Coimbra, C.F.M. Forecasting of global and direct solar irradiance using stochastic learning methods, ground experiments and the NWS database. *Sol. Energy* **2011**, *85*, 746–756.
19. Martin, L.; Zarzalejo, L.F.; Polo, J.; Navarro, A.; Marchante, R.; Cony, M. Prediction of global solar irradiance based on time series analysis: Application to solar thermal power plants energy production planning. *Sol. Energy* **2010**, *84*, 1772–1781.
20. Boata, R.S.; Gravila, P. Functional fuzzy approach for forecasting daily global solar irradiation. *Atmos. Res.* **2012**, *112*, 79–88.
21. Dambreville, R.; Blanc, P.; Chanussot, J.; Boldo, D. Very short term forecasting of the Global Horizontal Irradiance using a spatio-temporal autoregressive model. *Renew. Energy* **2014**, *72*, 291–300.
22. Diagne, M.; David, M.; Lauret, P.; Boland, J.; Schmutz, N. Review of solar irradiance forecasting methods and a proposition for small-scale insular grids. *Renew. Sustain. Energy Rev.* **2013**, *27*, 65–76.
23. Zeng, J.; Qiao, W. Short-term solar power prediction using a support vector machine. *Renew. Energy* **2013**, *52*, 118–127.
24. Dong, Z.; Yang, D.; Reindl, T.; Walsh, W.M. Short-term solar irradiance forecasting using exponential smoothing state space model. *Energy* **2013**, *55*, 1104–1113.
25. Ren, Y.; Suganthan, P.N.; Srikanth, N. A Comparative Study of Empirical Mode Decomposition-Based Short-Term Wind Speed Forecasting Methods. *IEEE Trans. Sustain. Energy* **2015**, *6*, 236–244.
26. Mellit, A.; Benghaneim, M.; Bendekhis, M. Artificial Neural Network Model for Prediction Solar Radiation Data: Application for Sizing Stand-Alone Photovoltaic Power System. In Proceedings of the 2005 IEEE Power Engineering Society General Meeting, San Francisco, CA, USA, 12–16 June 2005.
27. Farhad, S.G.; Khaze, S.R.; Malekl, I. A new Approach in Bloggers Classification with Hybrid of K-Nearest Neighbor and Artificial Neural Network Algorithms. *Indian J. Sci. Technol.* **2015**, *8*, 237–246.
28. Zhang, Y.; Wang, J. GEFCom2014 Probabilistic Solar Power Forecasting based on k-Nearest Neighbor and Kernel Density Estimator. In Proceedings of the 2015 IEEE Power & Energy Society General Meeting, Denver, CO, USA, 26–30 July 2015.
29. López, G.; Battles, F.J.; Tovar-Pescador, J. Selection of input parameters to model direct solar irradiance by using artificial neural networks. *Energy* **2005**, *30*, 1675–1684.
30. Anil, K.; Jain, J.; Mao, K. Mohiuddin, Artificial Neural Networks: A Tutorial. *Computer* **1996**, *29*, 31–44.
31. Wu, B.; Li, K.; Yang, M.; Lee, C-H. A Reverberation-Time-Aware Approach to Speech Dereverberation Based on Deep Neural Networks. *IEEE/ACM Trans. Audio Speech Lang. Process.* **2017**, *25*, 102–111.

