

## Data Analytics on Solar Energy Using Hadoop

K.Akila<sup>1</sup>, B.Manasavani<sup>2</sup>, Dr.Gopikrishnan M<sup>3</sup>

<sup>1</sup>Department Of ComputerEngineering Prathyusha Engineering College Thiruvallur, B.E.,India.

<sup>2</sup>Department Of ComputerEngineeringPrathyusha EngineeringCollege

**ABSTRACT:** Missing data is one of the major issues in data mining and pattern recognition. The knowledge contains in attributes with missing data values are important in improving regression correlation process of an organization. The learning process on each instance is necessary as it may contain some exceptional knowledge. There are various methods to handle missing data in regression correlation. Analysis of photovoltaic cell, Sunlight striking on different geographical location to know the defective or connectionless photovoltaic cell plate. we mainly aim to showcasing the energy produced at different geographical location and to find the defective plate. And also analysis the data sets of energy produced along with current weather in particular area to know the status of the photovoltaic plates. In this project, we used Hadoop Map-Reduce framework to analyze the solar energy datasets.

**Keywords:** missing data, genetic algorithm, regression correlation.

### INTRODUCTION

Missing data is the missing form of information about phenomena, which is important, and it is the information in which we are interested. The existence of missing data is one significant problem in data quality. Data quality plays a major role in machine learning, data mining and knowledge discovery from databases. Machine learning algorithms handle missing data in a quite naive way. To avoid biasing in induced hypothesis missing data treatment should be carefully handled. Imputation is a process that replaces the missing values in an instance by some reasonable values. The case substitution is the method developed for dealing with missing data in instances and it has some drawbacks when applied to the data mining processes. The methods, such as substitution of missing values by the attribute mean or mode should be cautiously handled to avoid inclusion of bias.

#### 1.1 Randomness of Missing Data

Missing data randomness is classified [1] in three classes.

**Missing completely at random (MCAR):** Missing values are scattered randomly across all instances. In this type of randomness, any missing data handling method can be applied without risk of introducing bias on the data. It occurs when the probability of an instance having a missing value for an attribute does not depend on either the known values or the missing data. Differences on attributes establish that the two groups do not differ significantly.

**Missing at random (MAR):** Missing at random (MAR) is a condition, which occurs when missing values are not randomly distributed across all observations but are randomly distributed within one or more classes (ex. missing more among whites than non-whites, but random within each). The probability of an instance with a missing value for an attribute may depend on the known values, and not on the value of the missing data itself.

**Not missing at random (NMAR):** Not missing at random is the most challenging form, occurs when missing values are not randomly distributed across observations. It is also called as non-ignorable missingness. The probability of an instance with a missing value for an attribute might depend on the value of that attribute.

#### 1.2 Handling Missing Data

Missing data handling methods are categorized as follows **Ignoring data:** This method throw-outs all instances with missing data. There are two core methods to discard data with missing values. The first one is known as complete case analysis. It is available in every one of statistical packages and is the default method in many programs.

The next method is recognized as discarding instances or attributes. This method determines the level of missing data on each instance and attribute, and deletes the instances or attributes with high extents of missing data.

Prior to deleting any attribute, it is vital to evaluate its connotation to the investigation. The methods, complete case analysis and discarding is executed only if missing data are missing completely at random. The missing data that are not missing completely at random contain non-random elements that may prejudice the results.

**Imputation:** In imputation-based procedures missing values are imputed with reasonable, probable values rather than being deleted totally. The objective is to use known associations that can be recognized in the valid range values of the dataset to facilitate estimating the missing values.

Data quality plays a major role in machine learning, data mining and knowledge discovery from databases. Machine learning algorithms handle missing data in a quite naive way. To avoid biasing in induced hypothesis missing data treatment should be carefully handled. Imputation is a process that replaces the missing values in instance by some reasonable values. In first method limitation of  $\Delta T$  was adjusted, above equation was used for  $\Delta T < 8^\circ\text{C}$ . Table 1 shows criteria for the decision of  $\tau$  value.  $\tau$  value of 0.6-0.7 are commonly used for clear sky atmospheric transmittance coefficient value. In this study  $\tau$  value of 0.69 was used for clear sky, assumed that the clear sky condition occurred when  $\text{RH} < 40\%$  and ambient temperature more than  $8^\circ\text{C}$ . Calculation algorithm was built based on decision matrix and the  $\tau$  value was locally determined using the training of data set to get minimum error.

The second method used in this study is by finding the correlation between RH, clearness index and beam transmittance. The data used to find correlation between beam transmittance and clearness index is measured data from new radiometer set that was installed in 2010 on the rooftop of Block P, Universiti Teknologi Petronas, which located about 30 km from Ipoh city. About 1 month, 5 minutes time step data of global, beam and diffuse radiation from June to July 2010 was used. Before find the correlation of beam transmittance and clearness index, RH-clearness index correlation was obtained from Ipoh city available data as can be seen in Figure 3.

Then beam transmittance-clearness index correlation can be obtained by scatter plot as can be seen in Figure 4. To plot Figure 4 some data were rejected due to obvious error that can be analyzed from measurement results, and the basic concept of terrestrial solar radiation characteristics. There are two core methods to discard data with missing values. The first one is known as complete case analysis. It is available in everyone of statistical packages and is the default method in many programs.

**Multiple Imputations:** It is a method by Rubin [1] for making multiple simulated values for each incomplete information, and then iteratively examining datasets with each simulated value substituted in every turn. The intention is, possibly, to generate estimates that better indicate true variance and uncertainty in the data than do regression methods. This permits expert staff and software to be used to create imputed datasets that can be analyzed by relatively naive users equipped with standard software. It can be very effective particularly for small to moderate level of missingness, where the missing data mechanism is organized, and for datasets that are to be placed in the public domain.

Not missing at random is the most challenging form, occurs when missing values are not randomly distributed across observations. It is also called as non-ignorable missingness. The probability of an instance with a missing value for an attribute might depend on the value of that attribute.

**Following constraint were used as data rejection criteria:**

- Reject night data
- Reject data if clearness index  $> 1$
- Reject data if beam transmittance  $> 1$
- Reject data if beam transmittance  $>$  clearness index
- Reject data if clearness index  $> 0.6$  and beam transmittance  $< 0.1$
- Reject data if clearness index  $< 0.2$  and beam transmittance  $> 0.15$

Correlation between RH and beam transmittance was obtained from above correlation and plotted in Figure 5. Balaras et al studied the relationship between beam transmittance and clearness index in Athens, Greece [10], the results of the study was adopted to carry out second method in this study. Regression results were presented as follows:

There are various methods to estimate solar radiation. Satisfactory result for hourly solar radiation estimation was obtained by using atmospheric transmittance model [1] while other authors have used diffuse fraction [2] and clearness index models [3]. Parametric or atmospheric transmittance model requires details atmospheric characteristic information [4]. This model gives high accuracy for clear sky/cloudless conditions, which is leading some author to use this model to evaluate the performance of an empirical model under cloudless conditions [5]. There are numerous authors

proposed this kind of model as mentioned in [6]. However, pure parametric model was not used in this study, since there is no detail atmospheric condition data for the site.

Meteorological parameters frequently used as predictors of atmospheric parameters since acquiring detail atmospheric conditions require advanced measurement. Meteorological parameters such as sunshineduration, cloud cover, ambient temperature, relative humidity, and precipitation data have been used to estimate atmospheric transmittance coefficient in parametric model. This kind of model is called meteorological model

There are various ways to estimate solar radiation on certain area on the earth. Ambient temperature based estimation is widely used since ambient temperature data are measured in many weather stations. In this study, missing data were estimated based on ambient temperature measurement and used measured RH data as atmospheric transmittance of the determination criteria. that estimate hourly solar radiation based on developed Campbell and Norman method [2], was adapted in this study.

Results of the new model then compared with the existing temperature-based solar radiation prediction model as follow:

### A. H-S model

Hargreaves and Samani [14] conducted an initial study on using  $T_{max}$  and  $T_{min}$  to estimate solar radiation by the following equation:

$$G_{Th} = K_r (T_{max} - T_{min})^{0.5} G_{oh} \quad (9)$$

$K_r$  is an empirical coefficient, which was recommended to be 0.16 for interior regions and 0.19 for coastal regions. In this study  $K_r$  was locally determined using training data set.

### B. H-S-A model

Annandale et al [15] modified H-S model by introducing correction factor as follow:

$$G_{Th} = K_r (1 + 2.7 \cdot 10^{-5} Z) (T_{max} - T_{min})^{0.5} G_{oh} \quad (10)$$

$Z$  is elevation in  $m$  and  $K_r$  was locally determined.

## V. STATISTICAL ANALYSIS FOR MODEL VALIDATION

Estimation results validated using statistical parameters. Pearson correlation coefficient was calculated as routine correlation indicator. Residual error was calculated using RMSE (Root Mean Square Error) and also presented in NRMSE (Normalized Root Mean Square Error) as follows:

$$RMSE = [\sum \{Y_c - Y_o\}^2 / n]^{0.5} \quad (11)$$

$$NRMSE = RMSE / y_{max} - y_{min} \quad (12)$$

where,  $Y_c$  is predicted variable  $Y_o$  is measured variable,  $n$  is number of data,  $y_{max}$  is maximum measured data  $y_{min}$  is minimum measured data.

As an addition, index of agreement was calculated using equation below:

$$d = 1 - [\sum (x_i - y_i)^2 / \sum (|x_i - \bar{x}_i| + |y_i - \bar{y}_i|)]^{0.5} \quad (13)$$

where,  $x_i$  is predicted variable  $y_i$  is measured variable,  $\bar{x}_i$  is averaged predicted variable and  $\bar{y}_i$  is averaged measured variable.

## VI. RESULTS AND DISCUSSIONS

Calculation have been carried out using methods 1 and 2, statistical calculation analysis also has been performed. Table 2 shows statistical analysis results of both methods and results of existing method (H-S and H-S-A method). The most satisfactory results were obtained using method 1. Figure 3 shows scatter plot of predicted and measured data for first method. The minimum RMSE value of 87.6 Watt/m<sup>2</sup> was obtained with 0.95 correlation coefficient and 0.97 index of agreement value. Previous method which use precipitation data obtained averaged index of agreement of 0.95, thus the model presented in this study also performed well.

**Table II. Statistical analysis results**

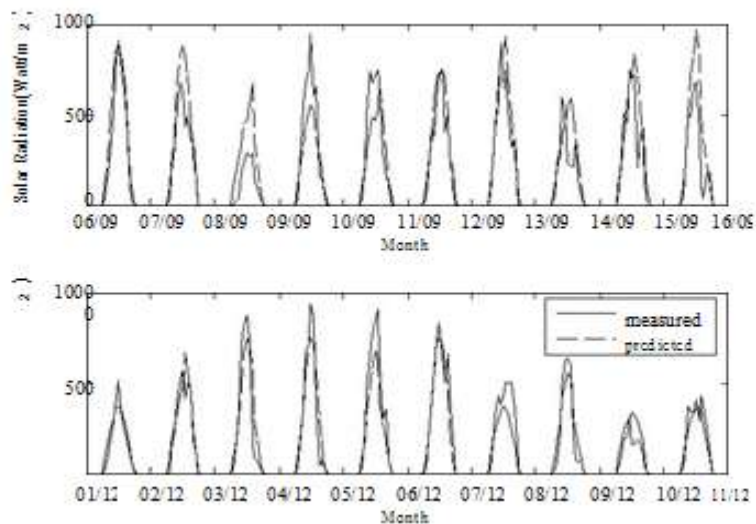
The methods are based on the decomposition model that is calculating each of solar radiation components, which depend on atmospheric transmittance. Two methods to predict atmospheric transmittance value using available meteorological data were proposed. In the first method, a decision matrix was used, while in the second method, regression correlation

one of meteorological parameters was used. The calculation results were evaluated using statistical parameter. Though the result shows both of the methods perform well, more satisfactory results were obtained from first method with Root Mean Square Error of  $87.6 \text{ Watt/m}^2$ , Normalized Root Mean Square Error of 8.29%, correlation coefficient of 0.95 and index of agreement of 0.97. Furthermore, the first method only needs ambient temperature and relative humidity data that common measured in meteorological stations. Graphical comparison of measured and predicted solar radiation for random dates. As can be seen most of the error occurred in low radiation or cloudy sky which often occurred in rainy season (Fig. 7-b). Then after the results has been validated using available data, the predicted results in the missing data days can be occupied in the measured data set. Now the time series composite data of solar radiation is completed and can be used for other purpose such as prediction of PV system performance for respected location.

Missing data is one of the major issues in data mining and pattern recognition. The knowledge contains in attributes with missing data values are important in improving decision-making process of an organization. The learning process on each instance is necessary as it may contains some exceptional knowledge. There are various methods to handle missing data in decision tree learning. Analysis of photovoltaic cell, Sunlight striking on different geographical location to know the defective or connectionless photovoltaic cell plate. we mainly aims To showcasing the energy produced at different geographical location and to find the defective plate.

And also analysis the data sets of energy produced along with current weather in particular area to know the status of the photovoltaic plates. In this project, we used Hadoop Map-Reduce framework to analyze the solar energy datasets.

Not missing at random is the most challenging form, occurs when missing values are not randomly distributed across observations. It is also called as non-ignorable missingness. The probability of an instance with a missing value for an attribute might depend on the value of that attribute



(b) Rainy season

**Figure 7:** Graphical comparison between measured and predicted solar radiation for random dates (Method 1)

Figure 8 shows prediction results in the days which solar radiation measurement were absent. In our case the amount of missing solar radiation data is 21 full days and 2 half days.

Data quality plays a major role in machine learning, data mining and knowledge discovery from databases. Machine learning algorithms handle missing data in a quite naive way. To avoid biasing in induced hypothesis missing data treatment should be carefully handled. Imputation is a process that replaces the missing values in an instance by some reasonable values. The case substitution is the method developed for dealing with missing data in instances and it is having some drawbacks when applied to the data mining processes. The method developed in this study may be applied for any location on the earth with notes, for first method the assignment criteria of atmospheric transmittance using RH and ambient temperature should be adjusted to the available solar radiation data of the area to get minimum error. To generate general criteria of atmospheric transmittance assignment using RH and ambient temperature further research is required with sufficient large amount of data for various areas. For the second method the correlation should be rebuilt based on available measurement nearest from the location to give satisfactory estimation results.

## VII. CONCLUSIONS

This paper represents to estimate solar radiation from available weather data. The methods presented in this paper are based on atmospheric transmittance determination using available meteorological data. In this method regression correlation of meteorological parameter was used. The key for the accuracy of above method is in the determination of beam atmospheric transmittance ( $\tau$ ). Beam atmospheric transmittance is the percentage of the beam (direct) radiation that will penetrate the atmosphere without being scattered. Kurt and Spokas used precipitation data to built decision matrix of atmospheric transmittance .temperature data. The second method is by using RH-clearness index, clearness index-beam atmospheric transmission and beam atmospheric transmission-RH correlation. The result shows that both methods perform well. Method 1 provided better results with minimum correlation coefficient of 0.95, RMSE of 87.6 Watt/m<sup>2</sup>, NRSME of 8.29% and index of agreement of 0.97. The prediction was intended to fill missing data in solar radiation data set to get complete time series data. However, in this study only one year of one area data have been used. Validation using sufficient large amount of data is required for wider application of the method.

## REFERENCES

- [1]. G. S. Campbell and J. M. Norman, "Introduction to Environmental Biophysics. 2nd ed. New York: Springer-Verlag. Pp. 167-183, 1998
- [2]. Guofeng Wu, et al. "Methods and strategy for modeling daily global solar radiation with measured meteorological data – A case study in Nanchang station, China", *Energy Conversion and Management* 48, 2447-2452, 2007
- [3]. Yang K, Koike T. Estimating surface solar radiation from upper-air humidity. *Solar Energy*, 72(2):177-86, 2002
- [4]. N. Mohan Kumar et al, "An empirical model for estimating hourly solar radiation over the Indian seas during summer monsoon season", *Indian Journal of Marine Sciences*, Vol. 30, pp 123-131, 2001
- [5]. Reindl DT, Beckman WA, Duffie JA, "Diffuse fraction correlations", *Solar Energy* 1990; 45:1-7
- [6]. F.J. Batlles et al, "Empirical modeling of hourly direct irradiance by means of hourly global irradiance", *Energy* 25: 675-688, 2000
- [7]. Kurt Spokas and Frank Forcella, "Estimating hourly incoming solar radiation from limited meteorological data", *Weed Science*, 54:182- 189, 2006
- [8]. Hunt LA, Kuchar L, Swanton CJ. Estimation of solar radiation for use in crop modelling. *Agric. Forest Meteorol.* 91(3-4):293-300, 1998
- [9]. Louche et al, "Correlations for direct normal and global horizontal irradiation on French Mediterranean site", *Solar Energy* 46: 261-266, 1991 as cited by [8]
- [10]. Athens, Greece. et al., "On the relationship of beam transmittance on clearness index for Athens, Greece", *Int. J. Solar Energy*, Vol. 7, pp 171-179, 1989