# IMPROVED NEURAL NETWORK PREDICTION PERFORMANCES OF ELECTRICITY DEMAND: MODIFYING INPUTS THROUGH CLUSTERING

K.A.D. Deshani[1], Liwan Liyanage Hansen[2], M.D.T. Attygalle[3], A. Karunaratne[4]

[1, 3, 4]Department of Statistics, University of Colombo, Colombo 03, Sri Lanka
[2] School of Computing, Engineering and Mathematics, University of Western Sydney, Australia
[1]deshani@stat.cmb.ac.lk, [3]dilhari@stat.cmb.ac.lk, [4]ak@stat.cmb.ac.lk
l.liyanage@uws.edu.au

## ABSTRACT

*Accurate prediction of electricity demand can bring extensive benefits to any country as the forecast values help the relevant authorities to take decisions regarding electricity generation, transmission and distribution much appropriately. The literature reveals that, when compared to conventional time series techniques, the improved artificial intelligent approaches provide better prediction accuracies. However, the accuracy of predictions using intelligent approaches like neural networks are strongly influenced by the correct selection of inputs and the number of neuro-forecasters used for prediction. This research shows how a cluster analysis performed to group similar day types, could contribute towards selecting a better set of neuro-forecasters in neural networks. Daily total electricity demands for five years were considered for the analysis and each date was assigned to one of the thirteen day-types, in a Sri Lankan context. As a stochastic trend could be seen over the years, prior to performing the k-means clustering, the trend was removed by taking the first difference of the series. Three different clusters were found using Silhouette plots, and thus three neuro-forecasters were used for predictions. This paper illustrates the proposed modified neural network procedure using electricity demand data.*

## KEYWORDS

*Clustering, Silhouette plots, Improve performance*

## 1. INTRODUCTION

Predicting the future electricity demand is an essential task for a country, as a huge amount of money could be saved by utilizing the available electricity generation options. In this scenario, increasing the accuracy of short-term predictions is very crucial, as decisions regarding the required load, has to be taken within a short period of time. Literature regarding short-term load forecasting consists of both conventional time series models and artificial intelligent approaches from different fields mostly in the field of engineering. To develop a dynamic forecasting system, intelligent approaches yields better results than usual conventional time series techniques as they could be adapted to suit novel conditions and handle more complex patterns in data. However, the

accuracy of predictions using intelligent approaches like neural networks are strongly influenced by the correct selection of inputs and the number of neuro-forecasters used for prediction. This paper presents a cluster analysis, performed to group similar day types with respect to electricity demand. Even though many external causes like metrological conditions such as temperature, rainfall, humidity, wind speed and cloud cover, economic and demographic factors influence the electricity demand, this paper has considered only a single input which is day type. The main focus has been given to this input as this paper attempts to illustrate how data mining techniques can be complimented by statistical techniques to make them more efficient.

A dataset consisting of daily total electricity demands in Sri Lanka was considered for the period of 01st January 2008 to 31st December 2012. Each day was assigned to one of the predefined thirteen categories, suitable to Sri Lanka.

## 2. LITERATURE REVIEW

Literature related to load forecasting reveal that higher prediction accuracies could be obtained when using intelligent techniques when compared to using conventional statistical techniques (Farahat & Talaat, 2012; Barzamini, Hajati, Gheisari, & Motamadinejad, 2012; Nagi, Yap, Tiong & Ahmed, 2008). Many researchers point out the importance of using intelligent techniques in situations where quick weather changes lead to fail accurate predictions. (Seetha & Saravanan, 2007; Senjyu, Takara, Uezato, & Funabashi, 2002; Barzamini et al., 2012). Some of those popular intelligent techniques used in the literature are neural networks, fuzzy inference systems, genetic algorithms and expert systems.

Many researches had used the effect of different day types to enhance the load predictions considering their own country's situations. The literature reveals that, Soared and Medeiros (2008) had incorporated the maximum number of day types to their model, as Sunday - Saturday, holiday, working day after holiday, working day before holiday, working day between a holiday and weekend, Saturday after a holiday, working only during the mornings, working only during the afternoons and Special holidays. Another research considers seven days of the week and bank holidays as day types, and a principal component analysis had been performed accordingly and a segmentation scheme based on the first principal direction had been used to cluster similar months (Cho, Goude, Brossat, & Yao, 2013). Unlike these approaches, (Barzamini et al., 2012) had divided the weekly days into four categories based on unique load lags and had incorporated to the model. Considering the above, this research considers thirteen day types, which can be considered as different in Sri Lankan context.

Even though thirteen day types are considered, including all these day types into the model will complicate the prediction process. As such, the 'day type' will be clustered into similar day types in order to avoid complexities in the computation operations and to reduce forecasting error when training the neural networks (Barzamini et al., 2012; Seetha & Saravanan, 2007). They have discussed how accurate predictions are made when the inputs are wisely chosen to be fed into the neural network having different neuro-load forecasters to train similar featured loads. Literature also shows that in some research, similar days had been clustered based on experience of the experts of electricity supplying companies rather than performing any statistical analysis (Cho et al., 2013). Moreover, to understand energy consumption patterns in industrial parks, a cascade application of a Self-Organizing Map and a clustering k-means algorithm had been performed by Hernandez, Baladron, Aguiar, Carro & Esguevillas (2012). Even though no study has considered performing a cluster analysis, this study focuses on a statistical analysis based on k-means clustering to complement the neural network approach.

## 3. METHODOLOGY

### 3.1. Unit Root Test

Two common trend removal or de-trending procedures are first differencing and time-trend regression. First differencing is appropriate for I(1) time series and time-trend regression is appropriate for trend stationary

I(0) time series. Unit root tests can be used to determine if trending data should be first differenced or regressed on deterministic functions of time to render the data stationary. In this paper a stochastic trend could be seen in order to test for it Augmented Dickey Fuller test was used.

The ADF test tests the null hypothesis that a time series $Y_t$ is I(1) against the alternative that it is I(0), assuming that the dynamics in the data have an ARMA structure. The ADF test is based on estimating the test regression

$$Y_t = \beta' D_t + \phi Y_{t-1} + \sum_{j=1}^{p} \psi_j \Delta Y_{t-j} + \varepsilon_t$$

where $D_t$ is a vector of deterministic terms (constant, trend etc.), p is the lagged difference terms, $\Delta Y_{t-j}$ are used to approximate the ARMA structure of the errors, and the value of p is set so that the error $\varepsilon_t$ is serially uncorrelated.

When the null hypothesis of having a unit root cannot be rejected, the series can be made stationary by taking the first difference. One should pay special attention not to take the first difference to make a series trend stationary when there is a deterministic trend as it will introduce a non-invertible moving average component.

### 3.2. K-Means clustering

K-means clustering is a partitioning method. It partitions data into k mutually exclusive clusters. Unlike hierarchical clustering, k-means clustering operates on actual observations (rather than the larger set of dissimilarity measures), and creates a single level of clusters. The distinctions mean that k-means clustering is often more suitable than hierarchical clustering for large amounts of data.

Each cluster in the partition is defined by its member objects and by its centroid, or center. The centroid for each cluster is the point to which the sum of distances from all objects in that cluster is minimized. kmeans computes cluster centroids differently for each distance measure, to minimize the sum with respect to the measure that you specify.

Distance measure: In this paper, 'kmeans' function in Matlab software has been used with 'city block' as the distance measure. Even though there are five distance measure options, only 'city block' and 'sqEuclidean' were suited for the data, and the results for both the cases were almost similar. Therefore, 'sum of absolute differences', that is the 'city block' distance measure was considered.

Determining the number of clusters: To get an idea of how well-separated the resulting clusters are silhouette plot can be used. The silhouette plot displays a measure of how close each point in one cluster is to points in the neighboring clusters. This measure ranges from +1, indicating points

that are very distant from neighboring clusters, through 0, indicating points that are not distinctly in one cluster or another, to -1, indicating points that are probably assigned to the wrong cluster.

Avoiding Local Minima: Like many other types of numerical minimizations, the solution that kmeans reaches often depends on the starting points. It is possible for kmeans to reach a local minimum, where reassigning any one point to a new cluster would increase the total sum of point-to-centroid distances, but where a better solution does exist. However using 'replicates' one can overcome that problem by taking the one with the lowest total sum of distances, over all replicates as the final answer. (The MathWorks)

## 4. ANALYSIS AND INTERPRETATION

### 4.1. Trend Removal Process

Figure 1 displays fluctuations of daily total electricity demand from January 2008 to December 2012. A stochastic trend could be seen over the years.
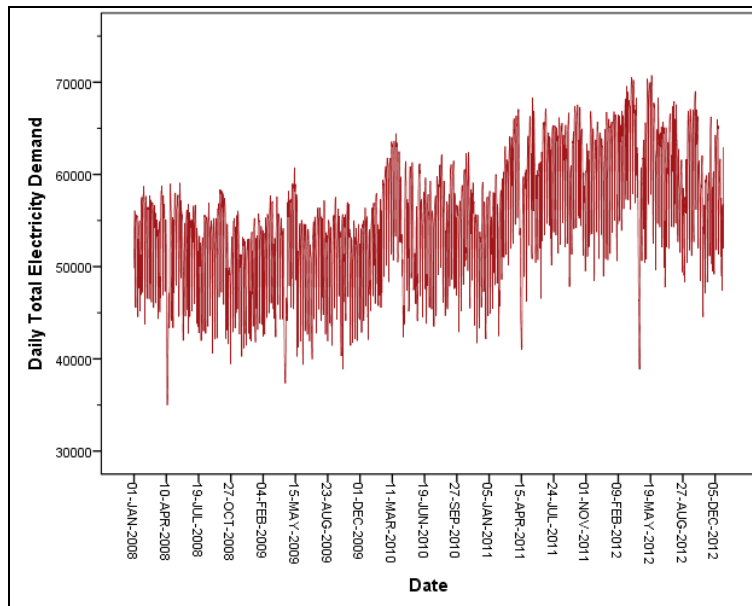


Figure 1.  Daily total electricity demand over the years

Figure 1 shows that the electricity demand gradually rises during this period. This was confirmed by the Augmented Dickey-Fuller Unit Root test. When the test was performed using the original series without any additional variables, the null hypothesis of having a unit root was not rejected suggesting that the original series is not stationary. However, when a trend and an intercept terms was included into the model, the null hypothesis of having a unit root was strongly rejected and both the trend and intercept coefficients were significant in the model. When applying the test for the first difference series, the null hypothesis of having a unit root was rejected even without the trend and the intercept term. Therefore, it can be concluded that there is a trend in the series and the trend can be removed using the first difference as a difference stationary process.

## 4.2. Clustering Similar Day Types

After taking the first difference of the time series, the trend was removed and then the days were categorized based on the first difference of the daily total electricity demand. The trend has been removed prior to clustering, as days in the latter years tend to cluster into separate clusters where the demand is comparatively higher than the former years.

In the dataset, each day had been assigned into one of the thirteen categories; Sunday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday, Poyaday, PBM Holiday (Public, Bank, Mercantile), PB Holiday (Public, Bank), Working day before holiday, Working day after holiday, Working day between a holiday and a weekend, Saturday after holiday.

In order to select the most appropriate number of clusters, silhouette plots were used based on the results of the k-means algorithm. Three clusters could be found as the most appropriate number of clusters, which resulted the maximum average Silhouette Value of 0.709384. (Figure 2). In order to avoid the iterations to end up at local minimas, each clustering procedure was replicated 5 times and considered the one with the lowest total sum of distances.
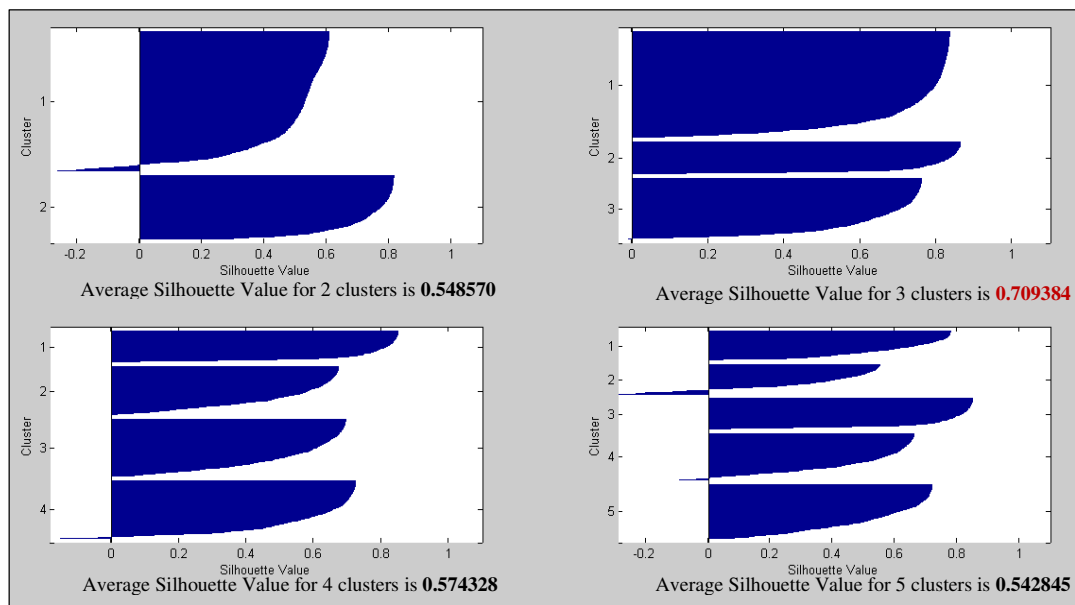


Figure 2.  Silhouette Plots to determine the correct number of clusters

Based on results of the cluster analysis, the thirteen day types could be categorized into three clusters (Table 1). Figure 3 displays how the data points are scattered across three layers. A very high percentage of Tuesdays (97.2%), Wednesdays (95.3%), Thursdays (95.8%), Fridays (96.7%) and days before holidays (92.1%) were clustered into the first cluster. A 98.2% of the Mondays were classified into the second cluster. Finally, a high percentage of Saturdays (91.1%), Sundays (91.4%), Poyadays (88.7%) and PBM holidays (74.4%) were clustered into the third cluster. However, PB holidays and Saturday after holidays seemed not prominent in any of these three clusters.

Table 1.  Distribution of day types across the three identified clusters

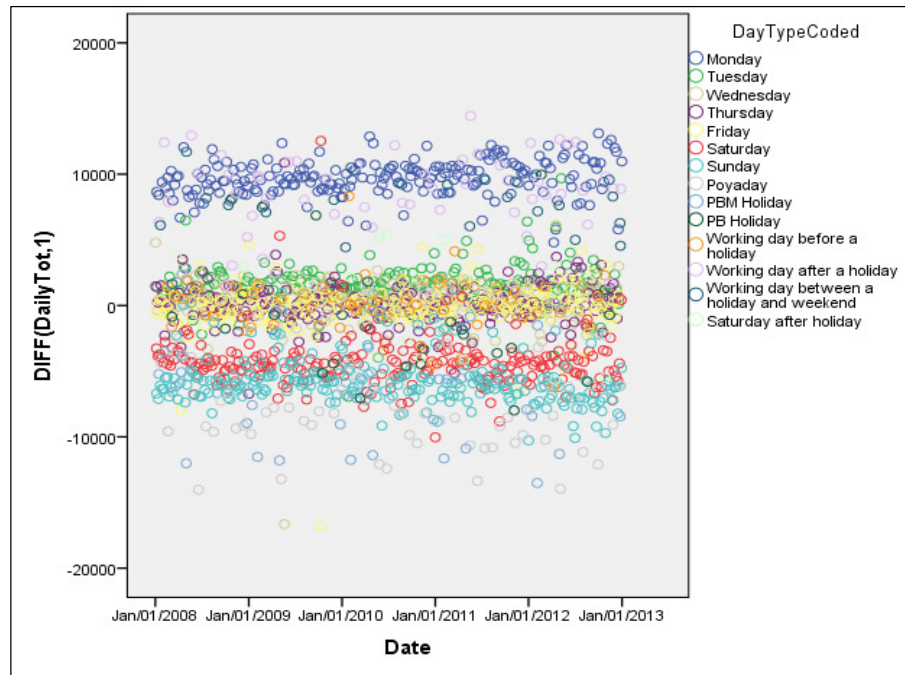| | | Clusters based on the first difference series | | |
| | | 1 | 2 | 3 |
| | | Row N % | Row N % | Row N % |
| Day Type | Monday | 1.8 | 98.2 | 0.0 |
| | Tuesday | 97.2 | 1.9 | 0.9 |
| | Wednesday | 95.3 | 0.5 | 4.3 |
| | Thursday | 95.8 | 0.5 | 3.8 |
| | Friday | 96.7 | 0.0 | 3.3 |
| | Saturday | 8.0 | 0.9 | 91.1 |
| | Sunday | 8.6 | 0.0 | 91.4 |
| | Poyaday | 11.3 | 0.0 | 88.7 |
| | PBM Holiday | 25.6 | 0.0 | 74.4 |
| | PB Holiday | 39.1 | 17.4 | 43.5 |
| | Working day before a holiday | 92.1 | 1.6 | 6.3 |
| | Working day after a holiday | 22.2 | 77.8 | 0.0 |
| | Working day between a holiday and weekend | 26.9 | 73.1 | 0.0 |
| | Saturday after holiday | 58.8 | 11.8 | 29.4 |



Figure 3. Spreading of different day types across the three clusters

**<u>Working day after a holiday</u> ( Cluster 1- 22.2%, Cluster 2 – 77.8%, Cluster 3 – 0.0%)**

There were 15 'day after holiday's classified into cluster 1 where normal Tues-Friday was included. The rest of the days were similar to Mondays contributing 77.8% of the total 'day after holidays' to cluster 2. It was interesting to find out that all these observations clustering into the

first cluster were either '*a day after a PB Holiday in a weekday*' or '*a working day after a new year day*'. Therefore, a new category was created as 'Workingday after a PB holiday in weekday or 'a working day after a new year day'.

Table 2.  Distribution of 'working day after a holiday's after modification

|  |  | Clusters based on the first difference series | | |
|---|---|---|---|---|
|  |  | 1 | 2 | 3 |
|  |  | Row N % | Row N % | Row N % |
| Day Type | Working day after a holiday I: (NOT PB Holidays) | 0.0 | 100.0 | 0.0 |
|  | Working day after a holiday II: (Working day after PB holiday in weekday OR Day after New Year) | 100.0 | 0.0 | 0.0 |

## Poya day ( Cluster 1- 11.3%, Cluster 2 – 0.0%, Cluster 3 – 88.7%)

It could be clearly seen that Poya days behaved like Saturdays and Sundays, except the 11.3% clustered under cluster 1 which behaved as normal working days like Tuesdays, Wenesdays, Thursdays and Fridays. From a detailed analysis of Poya days, it was found that, if a Poya day is a Monday there is a tendency to cluster those days into cluster 1. On the other hand, if a poya day is in May, it is the Wesak festival and if a poya day is in June, it is the Poson festival, which are exceptional to any other poya day. Therefore, a new category was introduced as 'Poyaday on Monday except in May and June'

Table 3.  Distribution of poya days after modification

|  |  | Clusters based on the first difference series | | |
|---|---|---|---|---|
|  |  | 1 | 2 | 3 |
|  |  | Row N % | Row N % | Row N % |
| Day Type | **Poyaday I :** Poyadays not in Mondays and Monday Poyadays in May & June | 1.8 | 0.0 | 98.2 |
|  | **Poyaday II:** Poya on Monday Except May & June | 100.0 | 0.0 | 0.0 |

## PBM Holiday ( Cluster 1- 25.6%, Cluster 2 – 0.0%, Cluster 3 – 74.4%)

When considering PBM holidays, all days following Wesak Full Moon poyadays had been classified to cluster I behaving differently than other PBM holidays. Therefore a new category was introduced as 'Day Following Wesak Full Moon Poyaday'.

Table 4.  Distribution of PBM holidays after modification

|  | | Clusters based on the first difference series | | |
|---|---|---|---|---|
|  | | 1 | 2 | 3 |
|  | | Row N % | Row N % | Row N % |
|  | PBM Holiday (Except day following Wesak) | 14.7 | 0.0 | 85.3 |
|  | Day following Wesak | 100.0 | 0.0 | 0.0 |

**Saturday after Holiday ( Cluster 1- 58.8%, Cluster 2 – 11.8%, Cluster 3 – 29.4%)**

This day type did not have prominent clustering percentages like other day types. Main problem was identified as the small number of data points. However, the percentages could be modified when the Saturday after holidays were classified into two new categories as 'Saturday after PB Holiday' and 'Saturday after PB Holiday OR Poyaday'.

Table 5.  Distribution of 'Saturday after holiday's after modification

|  | | Clusters based on the first difference series | | |
|---|---|---|---|---|
|  | | 1 | 2 | 3 |
|  | | Row N % | Row N % | Row N % |
|  | Saturday after PB holiday | 33.3 | 0.0 | 66.7 |
|  | Saturday after PBM holiday OR Poyaday | 66.7 | 25.0 | 8.3 |

**PB Holiday ( Cluster 1- 39.1%, Cluster 2 – 17.4%, Cluster 3 – 43.5%)**
PB holidays could not be further categorized as there were a small number of observations for the subcategories.
Newly identified day types and how they are distributed among the three clusters are presented in Table 6.

## 5. CONCLUSION

From the K-means clustering algorithm, day types are categorized into three clusters. Therefore three back propagated neuro-load forecasters are derived and used to train data of the three clusters to achieve higher performances. When predicting electricity demand values, identifying the appropriate day type based on table 6 and feeding it into the correct neuro-load forecaster is shown to yield better results. If the considered time duration could be expanded, results can be improved and will be more accurate as there will be more observations for the sub categories. Further, instead of selecting the day type based on the extracted day types, one can automate the classification of a new data point to one of the three predefined clusters in order to increase the accuracy.

Table 6. Distribution of new day types across the three identified clusters

| | | Clusters based on the first difference series | | |
| | | 1 | 2 | 3 |
| | | Row N % | Row N % | Row N % |
|---|---|---|---|---|
| Day Type | Monday | 1.8 | 98.2 | 0.0 |
| | Tuesday | 97.2 | 1.9 | 0.9 |
| | Wednesday | 95.3 | 0.5 | 4.3 |
| | Thursday | 95.8 | 0.5 | 3.8 |
| | Friday | 96.7 | 0.0 | 3.3 |
| | Saturday | 8.0 | 0.9 | 91.1 |
| | Sunday | 8.6 | 0.0 | 91.4 |
| | Poyaday I | 1.8 | 0.0 | 98.2 |
| | Poyaday II | 100.0 | 0.0 | 0.0 |
| | PBM Holiday (Except day following Wesak) | 14.7 | 0.0 | 85.3 |
| | PB Holiday | 39.1 | 17.4 | 43.5 |
| | Working day before a holiday | 92.1 | 1.6 | 6.3 |
| | Working day after a holiday I | 0.0 | 100.0 | 0.0 |
| | Working day after a holiday II | 100.0 | 0.0 | 0.0 |
| | Working day between a holiday and weekend | 26.9 | 73.1 | 0.0 |
| | Saturday after PB holiday | 33.3 | 0.0 | 66.7 |
| | Saturday after PBM holiday OR Poyaday | 66.7 | 25.0 | 8.3 |
| | Day following Wesak poya day | 100.0 | 0.0 | 0.0 |

# 6. ACKNOWLEDGEMENTS

# REFERENCES

[1]   Barzamini, R., Hajati, F., Gheisari, S., & Motamadinejad, M. B. (2012). Short Term Load Forecasting using Multi-layer Perception and Fuzzy Inference Syatems for Islamic Countries. Journal of Applied Sciences , pp40-47.

[2]   Farahat, M. A., & Talaat, M. (2012). Short-Term Load Forecasting Using Curve Fitting Prediction Optimized by Genetic Algorithms. International Journal of Energy Engineering , pp23-28.

[3]   Hernandez, L., Baladron, C., Aguiar, J. M., Carro, B., & Esguevillas, A. S. (2012). Classification and Clustering of Electricity Demand Patterns in Industrial Parks. Energies , pp5215-5227.

[4]   Nagi, J., Yap, K. S., Tiong, S. K., & Ahmed, S. K. (2008). Electrical Power Load Forecasting using Hybrid Self-Organizing Maps and Support Vector Machines. International Power Engineering and Optimization Conference, (pp. 51-56). Selangor.
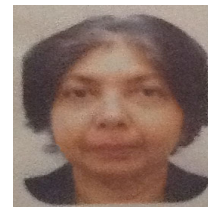
[5]    Seetha, H., & Saravanan, R. (2007). Short Term Electricity Load Prediction Using Fuzzy BP . Journal of Computing and Information Technology , pp267-282.

[6]    Soares, L. J., & Medeiros, M. C. (2008). Modeling and Forecasting short-term electricity load: A comparison of methods with an application to Brazilian data. International Journal of Forecasting , pp630-644.

[7]    The MathWorks, I. Statistics Toolbox.

## Authors

Ms. K.A.D. Deshani is a Lecturer (Probationary) in the Department of Statistics, University of Colombo since March 2011. Before she was absorbed to the permanent cadre, she has been working in the same department as an Assistant Lecture, obtaining a B.Sc Special Degree in Statistics with Computer Science in 2008. She has a keen interest in developing computer systems to incorporate the dynamic nature to statistical interpretations. Her research interests are in the areas of Operational Research and Data Mining. Currently, she is a member of the research team of the project titled "Developing an Economical Strategy for the Future Electricity Generation Procedure in Sri Lanka"; which received the University of Colombo research grants 2011 which is carried out in collaboration with University of Western Sydney, Australia and Ceylon Electricity Board, Sri Lanka. In June 2012 she registered for a M.Phil under the said grant. With her interest in research, she was a contributed speaker at the International Statistics Conference on the publication titled "Analysis of Efficiency of a Multi-Queue against a Single Queue with Many Servers: A Study on Advertisement Counter Queues at a Leading Newspaper Company" and was published in the Proceedings of the International Statistics Conference 2011, Colombo Sri Lanka. In 2013 a research paper titled "A Study of the Dynamic Behaviour of Daily Load Curve for Short Term Predictions" was published in the proceedings of the International Symposium for Next Generation Infrastructure (ISNGI) Australia.

Dr Liwan Liyanage joined University of Western Sydney in the year 1989 with university level teaching experience at University of Colombo, University of Wollongong and King Saud University Riyadh, totalling 12 years. Qualifications: B. Sc (First Class), Graduate Diploma in Applied Statistics, Masters Degree in Theoretical Statistics and the Ph. D. in the area of Applied Probability gives her the breadth of coverage across the statistics disciplines. At UWS she has been instrumental in developing many degree programs in particular the integrated degree B. Maths and IT using data mining as the integrating tool.  Senior lecturer (1995) and head of program of B. Maths & IT (1999). Her PhD was in random walk models, diffusion and related applications namely queuing theory and game theory. Thus her initial research was in applied probability, namely random walk models with difference equations, the master equation models with partial differential equations, and queuing models. This leads to differential equations representing diffusion and double diffusion. Her research, bridge the probabilistic models to the differential equation models of diffusion.  Her passion to integrate disciplines and research methods have led her to the current research areas which include innovative work in "Operational Statistics" a new area developed in collaboration with UC Berkeley; Optimisation Techniques and Data Mining. Application areas include bio security, public health, climate change, electricity production and demand. From her 10 PhD/Masters students 8 had completed the research successfully. She has established ongoing national and international linkages and research collaborations and 30+ publications and a book chapter. Her paper on Operational Statistics was the 2nd most downloaded paper in April 2006 from Science Direct's TOP 25 articles.

Dr. Attygalle has been Head of the Department of Statistics from 2010 to present. As a Senior Lecturer attached to the Department of Statistics from 2006 to present, she has been routinely involved in Teaching, Research and many other administrative roles such as being the Coordinator of the MSc in Applied Statistics, BSc Special degree and Joint Special degree programmes conducted by the Department of Statistics. She holds professional memberships of the Sri Lanka Association for the Advancement of Science (SLASS) and the Institute of Applied Statistics -Sri Lanka.

Dr. Attygalle obtained her PhD in Statistics from the Lancaster University, UK, and a MSc in Statistics from the Warwick University, UK, in 2006 and 1996 respectively. Prior to this she completed a Diploma in Applied Statistics, from the University of Colombo in 1992. She obtained her first degree majoring in Statistics, Applied Mathematics and Pure Mathematics also from the University of Colombo, graduating with a first class in 1987. As professional qualifications she has obtained the Staff and Educational Development Association (SEDA)-UK accreditation as a teacher in higher education in 2005 and the Certificate in Teaching in Higher Education (CTHE) by the Staff Development Centre of the University of Colombo also in 2005. Her key research areas are Statistical Modelling, Model Diagnostics, Data Visualization and Sports Statistics. As a senior lecturer she has supervised many undergraduate and postgraduate research projects. Currently she is so-supervising two MPhil/PhD research students. She had been instrumental in developing industry links with many private and government organisations over the years and thus has carried out many consultancy projects and other training programs through the Department of Statistics. She had also won one of the University of Colombo research grants in 2011.

Mrs. A. Karunarathne is a former head of the Department of Statistics, and also had served as the former head of Department of Statistics and Computer Science. She had been in the service   for more than 40 years and has been the key person to start the Special Degrees in Statistics and also initiate the Internship program  in the Department. As a senior lecturer she is conducting lectures mainly in the field of Operational Research and Stochastic Processes.

She obtained a Diploma di. Sp.(Operational Research) from University of Rome and her first degree was B.Sc.(Mathematics) from the University of Colombo. She has contributed to the continuous development and transmittance of statistical knowledge through many diverse avenues, a key example of which is her involvement in the publication of a book on basics of statistics titled "Moolika Sankayanaya" written for University entrants and A/L Science Students in Sept 1997. As a senior lecturer she has supervised many undergraduate and postgraduate research  projects.. Her key research areas are Stochastic Processes, Simulation, Queuing Models and Performance Modelling of Communication Networks.