

DEVELOPMENT OF A DATA-DRIVEN MODEL-BASED ANOMALY DETECTION
METHOD FOR WHOLE FACILITY LEVEL ENERGY USE DATA USING THE
ENERGY BALANCE LOAD VARIABLE

A Dissertation

by

HIROKO MASUDA

Submitted to the Office of Graduate and Professional Studies of
Texas A&M University
in partial fulfillment of the requirements for the degree of

DOCTOR OF PHILOSOPHY

Chair of Committee,	David E. Claridge
Committee Members,	Charles H. Culp
	Jeff S. Haberl
	Michael B. Pate
Head of Department,	Andreas A. Polycarpou

May 2018

Major Subject: Mechanical Engineering

Copyright 2018 Hiroko Masuda

ABSTRACT

The continuous monitoring of building energy use provides useful feedback to building owners and operators to achieve persistent energy efficiency and to make important engineering and financial decisions. The practice of managing the quality of metered data is essential to the successful utilization of energy data because metered energy data often contain errors and biases. This dissertation develops a method for automatically detecting anomalies in whole-building energy use data. The method can assist time-consuming and expensive tasks in monitoring and managing energy use data collected from a large portfolio of buildings. The method uses a variable called the energy balance load (E_{BL}), which is calculated from separately metered electricity, cooling, and heating energy use data. Anomalies are detected based on the distance between the E_{BL} value that is predicted by a data-driven reference model and the actual E_{BL} value. For the E_{BL} reference model, a simple regression model with weather variables is used so the model can be applied to various types of buildings with minimal information.

Updating reference models to account for the dynamic use and operations of buildings is a challenge. To address this challenge, this dissertation develops an anomaly detection method using the adaptive recursive least squares (RLS) filter as a reference model estimator and the standardized cumulative sum (CUSUM) test as a change detector. In the application to actual data, the new method demonstrates the ability to detect anomalies in a timely manner. The measurement bias in the chilled water use,

caused by a drift of the temperature sensor reading, was detected on the fourth day after the temperature began to drift. Furthermore, the start of a disabled time schedule for the heating, ventilation, and air conditioning (HVAC) systems was detected on the seventh day, and the change-back to the previous schedule was detected on the second day. Both the physical interpretation of the E_{BL} model parameters and the sensitivity and uncertainty analysis on the key parameters are presented such that they can be used as aids to warn analysts against physically impossible reference models.

ACKNOWLEDGEMENTS

I would like to thank my committee chair, Dr. David Claridge, and my committee members, Dr. Jeff Haberl, Dr. Charles Culp, and Dr. Michael Pate, for their guidance and support throughout the course of this research. I am particularly indebted to my advisor Dr. Claridge, who gave me freedom to explore various topics, understand problems, and develop solutions to the problems. His vision, guidance, and patience during my learning process greatly contributed to this dissertation.

This dissertation is based on the experience during the time I was working on the TAMU Project led by Dr. Juan-Carlos Baltazar at the Energy Systems Laboratory. I would like to gratefully thank Dr. Baltazar for all his help, for providing guidance, for his encouragement, and for making the research work possible. I am grateful to all of those with whom I have had the pleasure to work during this and other projects at the Energy Systems Laboratory. Completing this work would have been more difficult and less enjoyable were it not for the support and friendship provided by all of them.

Finally, I would like to thank Dr. Nobuo Nakahara, Dr. Harunori Yoshida, and other members of the Building Services and Commissioning Association, not only for supporting my decision to study at Texas A&M University, but also for providing opportunities to maintain professional relationships to the HVAC engineering community in Japan.

CONTRIBUTORS AND FUNDING SOURCES

Contributors

This work was supervised by a dissertation committee consisting of Dr. David E. Claridge (advisor), Dr. Charles H. Culp of the Department of Architecture, Dr. Jeff S. Haberl of the Department of Architecture, and Dr. Michael B. Pate of the Department of Mechanical Engineering. All work for the dissertation was completed independently by the student.

Funding Sources

There are no outside funding contributions to acknowledge related to the research and compilation of this document.

TABLE OF CONTENTS

	Page
ABSTRACT	ii
ACKNOWLEDGEMENTS	iv
CONTRIBUTORS AND FUNDING SOURCES.....	v
TABLE OF CONTENTS	vi
LIST OF FIGURES.....	ix
LIST OF TABLES	xiv
1. INTRODUCTION AND LITERATURE REVIEW.....	1
1.1 Introduction	1
1.2 Purpose	4
1.3 Objective	5
1.4 The Definition of the E_{BL}	5
1.4.1 Data Interval for E_{BL}	6
1.5 Literature Review.....	7
1.5.1 Thermal Energy Meters	7
1.5.2 Energy Anomaly Detection	9
1.5.3 Physical Significance of Regression Parameter Estimates for Building Energy Models	14
1.5.4 Energy-Balance-Based Variables	19
1.5.5 Previous Work on Data Screening Using the E_{BL}	21
1.5.6 Summary of Literature Review.....	26
1.6 Significance of the Study	29
2. STATISTICAL MODELING OF THE BUILDING ENERGY BALANCE VARIABLE	30
2.1 Introduction	30
2.2 Model Structure for the E_{BL}	31
2.2.1 Model Assumptions	31
2.2.2 Formulation of E_{BL} Model Structure	36
2.3 The Data	37
2.3.1 Description of the Data	37
2.3.2 Correlation Analysis of the Data	41
2.4 Regression Models	44

2.4.1	The Four-Parameter Change-Point (4P-CP) Model.....	44
2.4.2	Multiple Linear Regression (MLR) Models	46
2.4.3	Multiple Linear Models Incorporating an Autocorrelated Error Structure	47
2.4.4	Evaluation of the Fitted Models.....	49
2.5	Results and Discussion.....	51
2.5.1	Comparison of the Model Fits	52
2.5.2	Inclusion of AR(1) Error Structure in the MLR Model.....	57
2.5.3	Comparison of the Indoor Reference Temperature	59
2.5.4	Change Point Temperature	61
2.6	Chapter Summary.....	62
3.	ESTIMATION OF BUILDING PARAMETERS USING SIMPLIFIED ENERGY BALANCE MODEL AND METERED WHOLE-BUILDING ENERGY USE.....	65
3.1	Introduction	65
3.2	Formulation of Key Parameters	66
3.3	Data	69
3.3.1	Synthetic Data.....	69
3.3.2	Data from Actual Buildings.....	78
3.3.3	Estimation Procedure	83
3.4	Results and Discussion.....	85
3.4.1	Evaluation Using Synthetic Data.....	85
3.4.2	Application to the Data from Actual Buildings	91
3.5	Chapter Summary.....	95
4.	THE VARIATION OF THE INDOOR REFERENCE VARIABLE.....	97
4.1	Introduction	97
4.2	Methodology	98
4.2.1	Sensitivity and Uncertainty Analysis.....	98
4.2.2	Factorial Design.....	99
4.2.3	Building Energy Model	104
4.2.4	Model Selection	105
4.2.5	Analysis	106
4.3	Results and Discussion.....	108
5.	CONTINUOUS BUILDING ENERGY DATA MONITORING USING RECURSIVE LEAST SQUARES FILTER AND CUSUM CHANGE DETECTION: APPLICATION TO ENERGY BALANCE LOAD DATA.....	114
5.1	Introduction	115
5.2	Methodology	118
5.2.1	Energy Balance Load Model	118

5.2.2	Overview of Data Analysis Method	119
5.2.3	The RLS Filter	120
5.2.4	The CUSUM Test	127
5.3	Results and Discussion.....	130
5.3.1	General Performance of RLS Filters and Change Detection for EBL Data.....	130
5.3.2	Detection of Metering Error from Temperature Measurement Drift	135
5.3.3	Detection of HVAC Schedule Change	141
5.3.4	Limitations and Future Research	145
5.4	Chapter Summary.....	146
6.	CONCLUSIONS	148
6.1	Future Directions	152
6.2	Limitations	153
	REFERENCES	155
	APPENDIX A ENERGY USE AND E_{BL} DATA FOR REGRESSION ANALYSIS IN CHAPTER 2	169
	APPENDIX B ENERGYPLUS INPUT FILES FOR CHAPTER 3	173
	APPENDIX C ENERGY USE, E_{BL} AND Q_B DATA FOR REGRESSION ANALYSIS IN CHAPTER 3.....	174
	APPENDIX D FACTOR LEVELS AND E_{BL} PARAMETER ESTIMATES FOR SIMULATION RUNS IN CHAPTER 4	175
	APPENDIX E EQUEST INPUT FILES FOR PARAMETRIC SIMULATION RUNS IN CHAPTER 4.....	179
	APPENDIX F RECURSIVE LEAST SQUARES AND CUSUM CHANGE DETECTION PROGRAM IN CHAPTER 5	180

LIST OF FIGURES

	Page
Figure 1-1. Typical energy balance plot for data screening.....	24
Figure 1-2. Control limits of E_{BL} (Masuda et al. 2008).....	25
Figure 2-1. Semipartial correlation coefficients sr between the response variable E_{BL} and the explanatory variables T_{oa} , W_{oa}^+ , E_{sol} , and E_{ele} for the sample buildings. The buildings are ordered from the smallest to largest values for W_{oa}^+	43
Figure 2-2. The CV-RMSE values of the fitted models using data from the 56 sample buildings. A box represents the interquartile range (IQR) with the median as a horizontal line, and the whiskers extend from the ends of the box to the outermost data point that falls within the 75 th quantile +1.5·IQR and the 25 th quantile –1.5·IQR. The data outside this range are presented as outliers, and the building numbers are indicated. The broken line connects the means, and the inside error bars are constructed using 1 SE from the mean. The mean and SE values of each model are indicated in the plot.	53
Figure 2-3. The CV-RMSEs of the fitted models for the 10 buildings where the W_{oa}^+ variable had the greatest impact.	55
Figure 2-4. The CV-RMSEs of the fitted models for the 10 buildings where the W_{oa}^+ variable had the least impact.	55
Figure 2-5. E_{BL} residuals versus T_{oa} for 4P-CP(T), MLR(T,W), and AR1(T,W) models for building 31. The dashed lines represent approximate uncertainty intervals at $k = 2$ for the respective models.	56
Figure 2-6. The E_{BL} residuals versus T_{oa} for the 4P-CP(T), MLR(T,W), and AR1(T,W) models for building 25. The dashed lines represent approximate uncertainty intervals at $k = 2$ for the respective models.	57
Figure 2-7. Time series plots of the T_{oa} , the actual E_{BL} , the fitted E_{BL} by the MLR(T,W) model (CV-RMSE = 9.7%), and the fitted E_{BL} by the AR1(T,W) model (CV-RMSE =5.4%) for building 7.....	59
Figure 2-8. Time series plots of the T_{oa} , the actual E_{BL} , the fitted E_{BL} by the MLR(T,W) model (CV-RMSE=13.8%), and the fitted E_{BL} by the AR1(T,W) model (CV-RMSE =9.3%) for building 11.....	59

Figure 2-9. The T_{ref} of the fitted models using data from the 56 sample buildings. A box represents the interquartile range (IQR) with the median as a horizontal line, and the whiskers extend from the ends of the box to the outermost data point that falls within the 75 th quantile +1.5·IQR and the 25 th quantile –1.5·IQR. The data outside this range are presented as outliers, and the building numbers are indicated. The broken line connects the means, and the inside error bars are constructed using 1 SE from the mean. The mean and SE values of T_{ref} for each model are indicated in the plot.	60
Figure 2-10. The histogram and IQR of T_{cp} in the 4P-CP models for the 56 sample buildings plotted over the weather data: daily average W_{oa} versus daily average T_{oa}	61
Figure 3-1. System schedules for the As-is case.....	71
Figure 3-2. System schedules for the Ideal case	71
Figure 3-3. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the As-is case. The time series plots are in the top figure and scatter plots versus daily average outside air temperature are in the bottom figure.	72
Figure 3-4. E_{BL} and Q_B per unit conditioned floor area for the As-is case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL}	73
Figure 3-5. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the Ideal w/o solar case. The time series plots are in the top figure and scatter plots versus daily average outside air temperature are in the bottom figure.....	74
Figure 3-6. E_{BL} and Q_B per unit conditioned floor area in the Ideal w/o solar case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL}	75
Figure 3-7. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the Ideal w/ solar case. The time series plots are in the top and scatter plots versus daily average outside air temperature is in the bottom figure.	76
Figure 3-8. E_{BL} and Q_B per unit conditioned floor area in the Ideal w/ solar case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL}	77

Figure 3-9. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for Haas Hall. Time series plot is on the top and scatter plot versus daily average outside air temperature is on the bottom.	79
Figure 3-10. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for McFadden Hall. The time series plot is on the top and a scatter plot versus daily average outside air temperature is on the bottom.	80
Figure 3-11. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for Hobby Hall. Time series plot is on the top and a scatter plot versus daily average outside air temperature is on the bottom.	81
Figure 3-12. E_{BL} and Q_B daily data for Haas Hall during the period July 1, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL}	82
Figure 3-13. E_{BL} and Q_B daily data for McFadden Hall during the period September 1, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL}	82
Figure 3-14. The E_{BL} and Q_B daily data for Hobby Hall during the period July 21, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL}	83
Figure 3-15. Parameter estimates from synthetic daily data for the Ideal and As-is cases. For each of the parameters, the assumed true value is depicted as a solid line, and the parameter estimates using Q_B and E_{BL} are indicated as a circle and a cross respectively, along with the SEs, which are presented as bars.	87
Figure 3-16. Parameter estimates from synthetic monthly data for the Ideal and As-is cases. For each of the parameters, the assumed true value is depicted as a solid line, and the parameter estimates using Q_B and E_{BL} are indicated as a circle and a cross respectively, along with the SEs, which are displayed as bars.	88
Figure 3-17. The T_{ref} estimates and the distributions of T_{in} . For each case, estimates using E_{BL} and Q_B are depicted with the SEs, and the annual distribution of the daily average T_{in} is presented by box-and-whisker plots.	90

Figure 3-18. Daily average outside airflow rate (left) and T_{ref} (right) estimated for Haas Hall, comparing the estimates using daily and monthly interval data. Two different data periods are used. The SE is presented with bars for each estimate. 1 cfm = 1.699 m ³ /h.	93
Figure 3-19. Daily average outside airflow rate (left) and T_{ref} (right) estimated for McFadden Hall, comparing the estimates using daily and monthly interval data. The SE is depicted with bars for each estimate. 1 cfm = 1.699 m ³ /h.	93
Figure 3-20. Daily average outside air flow rate (left) and T_{ref} (right) estimated for Hobby Hall, comparing the estimates using daily and monthly interval data. The SE is displayed with bars for each estimate. 1 cfm = 1.699 m ³ /h.	94
Figure 4-1. Daily energy use and E_{BL} for Heep Laboratory between December 1, 2012 and December 31, 2013.	98
Figure 4-2. Impact of the level E_{BL} vs. T_{oa} stationary point on T_{ref}	111
Figure 5-1. Change-detection scheme. x = input, y = output, θ = unknown parameters, $\hat{\theta}$ = parameter estimates, and ε = residual.	120
Figure 5-2. The RMSEs for RLS filter using varied λ for three buildings. The minimum RMSE values are indicated by points, and dashed lines represent RMSEs for MLR models.	124
Figure 5-3. Flowchart for detecting changes in E_{BL} data	129
Figure 5-4. The result of RLS change detection ($\lambda = 0.9$, $h = 5$, $k = 0.5$) for CHA. The electricity (Elec), chilled water (Cool), and heating hot water (Heat) energy use and MLR model residuals are presented as well. Alarms are indicated as vertical lines.	134
Figure 5-5. Parameter estimates and standard prediction error estimates of the RLS filter for CHA. The regression estimates for the MLR model are presented as well.	135
Figure 5-6. Hourly average supply water temperatures of the chilled water energy meters for MDL and a neighboring building from April 1 to September 30, 2013. The difference of the MDL meter from the adjacent building meter is plotted in the top portion, and the values for individual meters are plotted in the bottom portion of the figure. Alarms are indicated on 23:00 for each of the alarmed days.	138

Figure 5-7. Energy use, RLS residuals, and CUSUM statistics for MDL during the period of April 1 to September 30, 2013 ($\lambda = 0.9$, $k = 0.5$, and $h = 5$).	139
Figure 5-8. Chilled water energy use for MDL from July 1 to July 31, 2013 is compared to the values during the previous year.	140
Figure 5-9. The E_{BL} for MDL from July 1 to July 31, 2013 is compared to the values during the previous year.	141
Figure 5-10. The E_{BL} for PVL from January 1 to December 31, 2013. Weekdays and weekends are plotted in different markers. The data between July 8 and September 2, 2013 are highlighted.	143
Figure 5-11. Energy use, RLS residuals, and CUSUM statistics of weekday and weekend processes for PVL during the period of December 1, 2012-December 31, 2013 ($\lambda = 0.9$, $k = 0.5$, and $h = 5$). The first month is the learning period.	144

LIST OF TABLES

	Page
Table 2-1. Gross floor area, number of daily E_{BL} data points, and functions of the sample buildings.	38
Table 2-2. Advantages and disadvantages of different models.	63
Table 3-1. Mathematical expressions for regression model parameters.	68
Table 3-2. Equations to calculate building parameters and the uncertainties from the regression estimates and SEs.	69
Table 3-3. Simulation input conditions for three synthetic datasets	71
Table 3-4. Measured outside airflow rates for three dormitory buildings along with the period and the number of days of available E_{BL} and Q_B data.	78
Table 3-5. Data sets used in the analysis, and the explanatory variable terms included in the regression models. The checked terms are included.	84
Table 3-6. Assumed true values and percent bias of estimates for m_v	88
Table 3-7. Assumed true values and percent bias of estimates for overall heat-loss coefficient U^*	88
Table 3-8. Assumed true values and percent bias of estimates for temperature slope $ \beta_T $	89
Table 3-9. The VIFs for explanatory variables in the models for synthetic data. The values for daily (D) and monthly (M) data are compared.	89
Table 3-10. Physical meaning of the reference parameter T_{ref} for different models used for synthetic data. Expected values are given as well.	91
Table 3-11. True values and bias of the m_v estimates for three dormitory buildings. The bias is expressed as a percentage of the measured value.	94
Table 3-12. The VIFs for explanatory variables in the models for three dormitory buildings. The values for daily (D) and monthly (M) data are compared.	94
Table 4-1. Factors and levels used in Shao (2006).	100
Table 4-2. Factors and levels used in Ji et al. (2008).	100

Table 4-3. Factors and levels (categorical factors).	102
Table 4-4. Factors and levels (continuous factors).	103
Table 4-5. The DSD factor assignment for each simulation run.....	103
Table 4-6. Normal distributions assumed for continuous factors. μ = mean, σ = standard deviation.....	107
Table 4-7. The final T_{ref} model and estimates selected based on the BIC. ‘×’ means interactions of two factors.	109
Table 4-8. The main effect and total effect indices for the final T_{ref} model. The bar chart illustrates the size of the total effect index.	110
Table 5-1. Forgetting factors having minimum RMSE (λ_{min}) for 15 selected buildings. The minimum RMSE from RLS (RMSE _{RLS}) and RMSE from MLR (RMSE _{MLR}) are compared.	124
Table A-1. Building number and filename.....	169
Table A-2. Data description	171
Table D-3. The levels of continuous factors in the simulation runs.....	175
Table D-4. The levels of categorical factors in the simulation runs.....	176
Table D-5. The parameter estimates of the E_{BL} regression model for each simulation run.....	177
Table E-6. Simulation run numbers and input files	179

1. INTRODUCTION AND LITERATURE REVIEW

1.1 Introduction

Growing environmental concern and increasing energy costs have been driving investments in building energy efficiency. In the US, the ESCO industry has grown since the 1990s, with an approximate annual growth of 9% for three years from 2009 to 2011 (Stuart et al. 2013); the green building market grew from 2% in 2005 to 44% in 2012 for commercial buildings (Fox and Morton 2013). Building energy codes are becoming more stringent: the 2012 International Energy Conservation Code (IECC) is estimated to result in 18.5% energy savings over the 2006 IECC on a weighted national basis for commercial buildings (Zhang et al. 2013). Energy savings from improved energy efficiency can pay back upfront costs within several years and provide financial benefits to building owners over the lifecycle of buildings. However, there is substantial evidence to suggest that buildings do not always yield their anticipated energy performance (Oates and Sullivan 2012; Newsham et al. 2009), and the effects of energy efficiency optimization and improvements such as retro-commissioning are not always persistent (Mills 2011). For these reasons, there is a growing interest in the continuous measurement and reporting of energy use to ensure that a building performs at the desired level. Federal energy policy (EPAct 2005; EISA 2007; EO 13514), building codes, standards, and voluntary certification programs, including the ASHRAE Standard 189.1: Standard for the Design of High-Performance Green Buildings (ASHRAE 2010a), the International Green Construction Code, and the Leadership in Energy &

Environmental Design (LEED) certification (USGBC 2014), have started to require the collection of energy use data on a continuous basis. The U.S. Department of Energy (DOE 2006) estimates energy savings from energy metering and the effective usage of collected data in large commercial buildings to be 1% to 20%.

Measured energy use data often contain anomalies and biases due to communication problems in the data acquisition network, the degradation or malfunction of metering and sensing devices, undesirable sensor locations, and incorrect factors used in the energy calculations. To reach useful conclusions from data, anomalies and biases must be flagged, and appropriate correction(s) must be made, if necessary. The periodic calibration of energy meters is an ideal practice to maintain metering accuracy; however, this is not always feasible due to the cost, especially when multiple buildings are centrally managed, as is the case with a university campus. In addition, continuous efforts to maintain the quality of a large volume of energy data can be time consuming and costly, creating a need for tools to automatically track the validity of the measured data and produce graphical and statistical summaries for review.

This dissertation develops a method to detect anomalies and biases that are worth investigating in whole-building energy data, and it demonstrates the application of this method to the energy data collected from a large number of buildings with minimal information. The method uses a variable called the energy balance load (E_{BL}) (Shao 2006; Shao and Claridge 2006), which is proposed for the initial stage of energy data screening. The E_{BL} variable is calculated from separately measured daily electricity, cooling, and heating energy use, and it principally represents the aggregate heat load in

the building. In the calculation of the E_{BL} variable, cooling and heating energy cancel out, which makes E_{BL} independent of air-side system types such as single duct, dual duct, constant air volume, and variable air volume.

To detect anomalies in energy data, Shao (2006) compared the E_{BL} values calculated from measured and simulated data. Using E_{BL} values instead of energy data directly allows one to use a simplified simulation that does not require detailed information of the air-side systems. However, when accurate building information is not available, the use of pre-tabulated values and other assumptions in the simulated load can lead to uncertainty, which decreases the ability to detect anomalies. Another approach proposed by Baltazar et al. (Baltazar et al. 2007) compared new E_{BL} data to the historical pattern of E_{BL} to detect anomalies, and it does not require building information. When the E_{BL} values for a building are plotted against the outside air temperature (T_{oa}), they demonstrate a largely linear pattern. In the plot, outliers and scattered patterns can be signs of possible metering problems. This data-driven approach using the E_{BL} variable was applied to the monthly validation process for metered energy consumption collected from Texas A&M University campus buildings for over 10 years, helping analysts to successfully detect various types of metering problems such as drift, calculation factor errors, and mislabeled sensors with minimal building system information (Baltazar et al. 2007; 2012).

Despite the usefulness of the graphical tool, the visual detection of anomalies by human analysts is a time-consuming task when hundreds of energy meters need to be verified continuously. In many cases of continuous energy metering in the building

industry, there is not sufficient resource allocation for data analysts to review data; therefore, there is a significant need for an inexpensive and reliable process to detect anomalies in energy data. This dissertation seeks to develop an automatable method of anomaly detections for energy use data collected from a large number of buildings with minimal building information, using the E_{BL} variable with a data-driven approach.

1.2 Purpose

The visual data screening that analysts perform using the E_{BL} graphical tool can be interpreted as a model-based detection (Isermann 1997), which is a change-detection algorithm that compares the model's prediction against the actual value and detects a change based on the size of the residuals. In the visual data screening using E_{BL} , an analyst visually compares the distance of new E_{BL} data from the recognized E_{BL} versus T_{oa} pattern to detect anomalies. Therefore, the E_{BL} versus T_{oa} pattern of historical data that the analyst recognizes in the plot is the prediction model. To replicate analysts' visual detection activity in an automatable process, this dissertation seeks to develop a method to estimate data-driven E_{BL} models and to detect sizable residuals.

In the data-driven approach, the E_{BL} model relies on past data; however, these past data are not always credible. Without continuous monitoring of the data quality, the past data can have biases and errors. Therefore, it is useful to have measures indicating physically unrealistic E_{BL} models. In the current usage of the E_{BL} versus T_{oa} plot (Baltazar et al. 2007, 2009), an analyst observes three characteristics of the pattern: the data concentration, the slope, and the T_{oa} at $E_{BL} = 0$, and sometimes compares the characteristics between different buildings to identify physically unrealistic patterns.

However, there has not been a study to quantify these characteristics of E_{BL} for standardized interpretation. Therefore, this dissertation will attempt to quantify and interpret the parameters that characterize the E_{BL} models.

1.3 Objective

The objective of this dissertation is to develop a model-based method for detecting anomalies in energy consumption data that are collected from a large number of buildings with minimal information, using the E_{BL} variable. This requires the following three steps:

- (a) develop a new data-driven modeling method for the E_{BL} variable that can be applied to different types of buildings to numerically describe E_{BL} patterns,
- (b) investigate the physical significance of the E_{BL} model parameter estimates, and
- (c) develop a new method to detect energy use anomalies using the E_{BL} models, and demonstrate the application.

1.4 The Definition of the E_{BL}

The definition of the E_{BL} variable is derived from the thermal energy flow through the whole building boundary. The net change of the total energy in a building, ΔE_{CV} , is equal to the difference between the total energy entering ($\Delta E_{entering}$) and leaving ($\Delta E_{leaving}$) the building. That is,

$$\begin{aligned}\Delta E_{CV} &= \Delta E_{entering} - \Delta E_{leaving} \\ &= Q_{air} + Q_{cond} + Q_{sol} + Q_{occ} + Q_{ele} - E_{cool} + E_{heat}\end{aligned}\tag{1.1}$$

where Q_{air} , Q_{cond} , Q_{sol} , Q_{occ} , and Q_{ele} are the building heat load components from air exchange, conduction through exterior surfaces, solar irradiance, and occupants, and thermal load from electrical energy use in the building respectively, and E_{cool} and E_{heat}

are the heat removed by cooling and added by heating respectively. When the time scale under study is long enough to diminish thermal inertia effects and to average out the change of the indoor air thermal condition, the system can be considered as quasi-steady state, and the left-hand side of Equation (1.1) yields 0. The E_{BL} is defined using the measurable variables as follows (Shao 2006):

$$\begin{aligned} E_{BL} &= E_{ele} - E_{cool} + E_{heat} \\ &= -Q_{air} - Q_{cond} - Q_{sol} - Q_{occ} \end{aligned} \quad (1.2)$$

where E_{ele} is the metered whole-building, non-cooling electricity use. For the evaluation of the E_{BL} , daily or longer interval energy use data are utilized to satisfy the quasi-steady state requirement by minimizing the thermal mass effects.

1.4.1 Data Interval for E_{BL}

The E_{BL} is calculated from daily or longer interval data, based on the assumption of a pseudo-steady state. The role of transients can be minimized based on the assumption that the initial and final conditions of daily interval energy data are nearly the same (Rabl 1988). Therefore, static and linear models can generally be sufficient for daily data. Hammersten (1984) estimated dynamic and static energy balance models using the data collected from existing houses, and he concluded that the use of static models is sufficient for time units of days or longer if the model is used for the prediction of energy consumption. Using measured hourly energy use data for 1 month, Kissock (1993) estimated the time lag between the net cooling load and the T_{oa} in a large institutional building to be 45 minutes, although the lighting load time constant is 5-10 hours for a typical commercial building. Based on this result, Kissock (1993) concluded

that the lag in the aggregate load is small for a large commercial building where the ventilation load is significant. Daily energy data still contain some dynamics resulting from changes in the HVAC controls and operation, building schedules, and thermal inertia. Weekly and monthly averaging can further reduce the dynamics; however, every averaging of data will result in a loss of information (Hammarsten 1987). In the following section, hourly and sub-hourly dynamic models are briefly reviewed; however, the main focus is on static models for daily or longer interval data. Based on the results of these earlier studies, steady-state models are used for daily E_{BL} data.

1.5 Literature Review

This section consists of six sub-sections. The first sub-section presents an overview of errors in thermal energy meters to highlight the importance of detecting anomalies in metered energy data. The next sub-section discusses the studies on energy anomaly detection, with an emphasis on the data-driven prediction models used in the model-based change-detection schemes. The third and fourth sub-section review studies on the use of the parameter estimates of data-driven energy models and those using variables that resemble E_{BL} respectively, while the fifth sub-section summarizes the previous studies on E_{BL} . The final sub-section is a summary of the literature review, and it describes new approaches investigated in this dissertation.

1.5.1 Thermal Energy Meters

Heat meters or thermal energy meters, such as chilled water and hot water meters, are known to be error-prone because the errors involving the temperature differential and the flow rate are exacerbated in the multiplication to calculate the heat

flow. A series of studies on the Measurement and Verification (M&V) process for a state-wide retrofit program called the LoanSTAR Program (O'Neal et al. 1990; Robinson 1992; Watt and Haberl 1994) found that the accuracy of thermal energy measurements was often compromised by inaccuracies in the flow measurements: errors in the manufacturers' pulse-per-gallon constants, errors due to an improper insertion depth of sensors, and errors due to a drop-out in the meter signal at low-velocity fluid flows. The flow meters installed in these studies were magnetic-type, tangential, paddlewheel flow meters. Watt and Haberl (1994) concluded that the magnetic-type tangential paddlewheel flow meters that were well-shielded and filtered for noise reduction were inaccurate at flow velocities lower than 2 feet per second.

Electromagnetic flow meters have increased the market share for the past decade (Zoebelein 2014; Choi et al. 2011), and they are widely used in thermal energy metering for commercial buildings today. The ease of maintenance due to a lack of in-flow parts and high accuracy are considered to be the strength of electromagnetic flow meters (Morris and Langari 2011; CSU 2012). Choi et al. (2011) tested 24 flow meters for hot water thermal energy meters, including turbine, electromagnetic, and ultrasonic types. They found high accuracy for electromagnetic meters; however, in their field tests, the electromagnetic type meter was inaccurate during the warmer season when the flow rate decreased below 6.9% of the maximum flow rate. The largest deviations for the 150 mm, 80 mm, and 50 mm diameter meters were approximately -70%, -60%, and -30% respectively. The authors recommended selecting a smaller diameter of electromagnetic flow meter to maintain better accuracy.

The temperature differential and the computation of the energy content in the chilled water or hot water flow are other major sources of inaccuracy in thermal energy meters. Typical temperature differential design conditions for the chilled water systems in commercial buildings are around 10°F (5°C). A differential as low as 5°F, or even smaller, is not unusual in actual buildings; therefore, temperature measurement accuracy is especially crucial to chilled water applications, since the temperature differential is small (CSU 2012; ASHRAE 2014). In an M&V case reported by Erpelding (2008), it was found that uncalibrated temperature sensors in the existing building energy management system (EMS) were mismatched by an average of 1.8°F (1°C), which caused a 19.7% average error in the chiller load measurement. CSU (2012) recommended using a BTU meter with factory matched integrated circuit temperature sensors that are extremely repeatable and linear, and that do not require calibration.

1.5.2 Energy Anomaly Detection

The detection of anomalies in energy consumption data was studied as a part of Fault Detection and Diagnostics (FDD) for building energy systems. Substantial research progress was made in the major research efforts sponsored by the International Energy Agency: Annex 25 (Hyvarinen 1996; 1997; Liddament 1999) and Annex 34 (Jagpal 2006), in which a number of analytical tools were developed to find problems in HVAC systems. The fault detection tools for building energy systems are classified based on their use of top-down and bottom-up approaches (Hyvarinen 1996; Friedman and Piette 2001; Ulickey et al. 2010). The top-down approach, sometimes referred to as energy

anomaly detection (Effinger et al. 2010), relies on whole-building energy use data to identify changes that resulted from various causes, whereas the bottom-up approach relies on BAS trend data to find system-level issues. This dissertation addresses the top-down approach, or energy anomaly detection, to screen unusual increases and decreases that are worth investigating. While there are numerous studies on the bottom-up approach to fault detection, they are not included in this review because the focus is on analyzing whole-building energy use data. Katipamula and Brambley (2005) and Kim and Katipamula (2017) provide a comprehensive summary of the bottom-up or system-level FDD approach.

The majority of energy anomaly detection methods found in the literature use the model-based approach. The model-based fault detection (Isermann 1997), which is a change-detection algorithm that commonly compares model predictions to actual values and detects a change based on the size of the residuals. This change-detection scheme is also known as analytical redundancy (Clark et al. 1975; Chow and Willsky 1984). The prediction model can be characterized as either physical or empirical. Physical models are based on physical principles with known parameters. In contrast, the parameters in empirical models are estimated using statistical methods based on input and output data. ASHRAE (2017) calls the physical modeling approach the forward (classical) approach and the empirical modeling approach the data-driven (inverse) approach.

Forward models calculate energy use based on the physical principles of building systems with known parameters. Maile et al. (2012) and O'Neill et al. (2014) utilize the measured data from a BAS as input into calibrated EnergyPlus simulation models, so the

models can emulate actual energy use behavior in semi-real time once the process is automated. While these tools can provide useful information for improvements and problem isolation, the modeling and calibration of detailed simulations require high-level simulation skills, along with comprehensive knowledge of buildings, to achieve accurate results. A fault-detection tool called ABCAT (Curtin 2007; Bynum et al. 2012) provides a low-cost method that can be implemented in existing buildings using limited measurements and simpler steady-state simulation at the trade-off of diagnostic granularity. The applications of energy anomaly detection using forward models found in the literature are not intended for a large number of buildings because of the requirement of accurate building information and calibration effort.

Unlike forward models, the data-driven models do not require prior knowledge of buildings because the unknown parameters are estimated from measured data. Therefore, the use of data-driven models is suitable for analyzing multiple buildings without detailed simulations; for example, applications to a portfolio consisting of multiple buildings. Various techniques are used to estimate data-driven models, depending on the granularity of the data and the complexity of the models. For daily or longer interval data, the regression technique is widely used. In the reviewed literature, Haberl and Claridge (1987), Haberl et al. (1988), Harris (1989), Stuart et al. (2007), and Liu et al. (2011) used regression models for energy anomaly detections.

Haberl and Claridge (1987) and Haberl et al. (1988) developed a tool to detect abnormal energy consumption in relation to previous performance, and they applied it to multiple institutional buildings. Their rule-based anomaly detection algorithm included

the distance of measured daily consumption from the predicted value using linear regression models for normal operation. Various building-specific calendar and schedule variables and weather data were included in the model. Haberl et al. (1989) later introduced the steady-state, change-point model as a function of ambient temperature to this fault-detection scheme, achieving better prediction of cooling and heating energy consumption.

Harris (1989) presented guides to monitor the energy performance of plants and buildings using the cumulative sum (CUSUM) (Page 1,954) of the difference between the measured values and the predicted values obtained through performance lines estimated with linear regression or other appropriate models. This model-based method using the residuals' CUSUM is called energy monitoring and targeting (M&T) or monitoring, targeting, and reporting (MT&R), and it is widely accepted as a monitoring method in the energy management field (Capehart et al. 2012). Stuart et al. (2007) analyzed half-hourly electricity consumption collected from 37 schools using the M&T approach, and they indicated that their method can detect changes that would be overlooked in a simple visual inspection of large datasets. The authors found that while CUSUM charts can reveal small-level shifts over time, there are some challenges in the interpretation of the CUSUM trend variation caused by skewed prediction.

Liu et al. (2011) developed a statistical toolkit to analyze monthly energy consumption for a portfolio of K-12 public school buildings in New York City. In this tool, the variable-based, degree-days (VBDD) regression models are estimated for each building, and the residuals are used for anomaly detection. The monthly seasonal

patterns are removed from the VBDD regression model residuals, and these residuals, after removing the seasonal patterns, are further fit to an autoregressive integrated moving average (ARIMA) model. The 95% confidence intervals of this ARIMA model for the residuals are used in conjunction with the predicted values of energy use as control limits to detect anomalies.

To predict hourly or shorter interval energy data, more complex modeling methods are often used to describe dynamics. Examples of such modeling techniques include the Artificial Neural Networks (ANN) (Yalcintas and Akkurt 2005; Yang et al. 2005; Dhar 1995; Karatasou et al. 2006), the Support Vector Machines (SVM) (Dong et al. 2005; de Wilde et al. 2013), the Seasonal Auto-Regressive Integrated Moving Average (SARIMA) (Henze et al. 2004), and the Fourier series (Dhar 1995; Dhar et al. 1999). There are some published applications of hourly or sub-hourly interval models for energy anomaly detections; for example, Dodier and Kreider (1999) used an Energy Consumption Index (ECI)—the ratio of actual to expected energy consumption as determined from a neural network—to detect whole-building energy anomalies. Yu and van Paassen (2003) used the residuals of predicted gas energy use with the Fuzzy Neural Network (FNN) models to detect a fault due to an open window in a room. Several studies indicate that strategically constructed linear regression models can predict hourly energy consumption as accurately as computationally more complex models (Kreider and Haberl 1994; Haberl and Thamilsaran 1996; Ramanathan et al. 1997). For the ease of computation and interpretation, Mathieu et al. (2011) chose the linear-regression-based energy prediction method over black-box models such as non-linear and time-

series models to predict 15-minute interval loads for their demand response evaluation tool, which is intended for large-scale application.

1.5.3 Physical Significance of Regression Parameter Estimates for Building Energy Models

The parameter estimates of data-driven building energy models can be physically explained with admittances such as heat-loss coefficients and solar apertures, time constants, and heat capacity. Rabl (1988) demonstrated that the parameter estimates for various dynamic, data-driven models, including thermal networks, modal analysis, autoregressive models, and the Fourier series, can be physically explained with admittances, time constants, and heat capacity, and these methods are basically equivalent to the problem of identifying the coefficients of a differential equation describing thermal balance. Integrating the differential equation results in a steady-state form with the parameters of heat loss and solar aperture coefficients (Hammarsten 1984; Rabl 1988). The E_{BL} is based on the assumption of the quasi steady-state, and the review of this section is limited to steady-state models.

The identification of building energy models was used to evaluate residential building envelope performance using short-time measurements. Several test procedures were developed, such as those by Palmiter et al. (1979), STEM (Subbarao 1988), Somogyi (1998), and Richalet et al. (2001). Many of the tests used some form of the co-heating test that Sonderegger and Modera (1979) developed to identify the heat-loss coefficient. Several studies (Bauwens et al. 2012; Butler and Dengel 2013; Bauwens and

Roels 2014) recently investigated the co-heating test using simplified quasi-stationary thermal models that utilized daily data. The co-heating tests involve heating the inside of a building to a constant, elevated temperature using electrical resistance heaters over a period of one to three weeks. Based on the daily average measurements of temperature differences between the inside and outside, global solar radiation, and energy supplied to the building, the overall heat-loss coefficient and global solar aperture coefficients were estimated using a linear regression analysis. The Building Research Establishment (BRE) in the UK (Butler and Dengel 2013) reported that the largest difference in the heat-loss coefficient estimate, from the value based on the known parameters, was approximately 17%; this statistic was a result of the tests that the BRE and six other teams conducted using test protocols based on the Leeds Metropolitan Protocol (Wingfield et al. 2010) at a pair of identical detached test houses. Adding aluminum solar shading to the windows appeared to reduce the difference in the heat-loss coefficient estimate to within -3.8% of the theoretical values, although this is not conclusive because the results are from only three tests.

Using utility data or energy consumption data collected through Building Energy Management Systems (BEMS), the Princeton Scorekeeping Method (PRISM), which is also an estimation method, was developed by Fels (1986) to estimate the variable-based degree-day (VBDD) models. Rabl and Rialhe (1992) incorporated occupancy variables into the VBDD model and applied it to 50 commercial buildings in five cities of France using the PRISM estimation. In many cases, the heat loss coefficients appeared to be under-estimated compared to theoretical values for many cases, indicating the heating

systems' efficiency exceeded 100% if the estimates and theoretical values are correct. The authors pointed out that excessive energy consumption due to the opening of windows during spring and autumn might be a possible reason. In the same study, the estimated base temperatures presented reasonable consistency with the thermostat set points. The authors concluded that the estimated models are reliable for energy prediction, and an interpretation of individual parameters can be a valuable tool; however, caution is advised.

Krarti (2012) suggested using the slope estimates of VBDD models as the overall heat-loss coefficient that takes into account any thermal coupling between heat conduction and air flow within the building envelope components for energy audits. Solupe and Krarti (2014) used the difference between the conventional and the regression heat-loss coefficient values as an estimate of the infiltration recovery factor. Here, the conventional value is the sum of the UA values for all exterior surface and infiltration loads, based on the blower door testing results. A similar use of the regression estimates of the heat-loss coefficient is found in Lowe et al. (2007), who analyzed the heat loss from air movement through cavities in party walls in masonry construction. The heat-loss coefficient, which was estimated from the co-heating test results from two houses, was a respective 75% and 103% larger than the calculated values based on the model that did not include cavities in party walls. This result was used as evidence of the significant heat loss through the cavity, which was not taken into account in conventional load calculation. These examples of using the values estimated by regression models demonstrate the possibility of identifying difficult-to-measure

physical building parameters. However, some authors question the accuracy of regression-estimated building parameters.

Hammarsten (1987) discussed the limitations of using steady-state energy regression models for building energy predictions and parameter estimations. The author points out an important drawback, namely that the physical interpretation of parameter estimates is only meaningful if the parameters are unbiased; however, there is a considerable risk of bias errors when simple energy balance models are used. Bauwens and Roels (2014) demonstrated that the solar aperture parameter loses much of its physical relevance in the simplification for regression models, and they advise that the solar aperture coefficient should not be calculated based on geometric and physical assumptions because the phenomena involved are complex and inextricably lumped into this parameter.

Some studies used the parameter estimates for simplified energy regression models as indices of energy performance benchmarking. These studies focused on obtaining crude information, rather than the physical accuracy of parameter estimates, regarding building energy performance from a portfolio of buildings. Casey et al. (2010) used the heating slope of the VBDD model as a metric to rank residential energy performance. Synthetic values of monthly natural gas consumption were generated for 567 homes using the DOE2 simulation program, which is designed for various home and occupant archetypes, and each simulation was labeled as compliant or non-compliant with the 2009 IECC. The VBDD model was estimated for each of the 567 homes, and the parameter estimates were compared. Although the heating slope estimated was

generally larger than the values calculated from the simulation inputs, the heating slope ranking identified 100% of the non-compliant homes; this is more effective than the 90% identification achieved using the annual fuel consumption. Liu et al. (Liu et al. 2011) used parameter estimates such as base load, balance-point temperature, and heating and cooling slopes for cooling and heating VBDD models to benchmark the relative building energy performance of 1,400 K-12 public schools. Their method analyzed the parameter estimates for different buildings using a multivariate analysis, with variables representing the building characteristics and operational activities. The stepwise variable selection procedure suggested that the following are related to electricity use: the gross floor area, the percentage of the area air conditioned, the number of students, the number of personal computers, the number of floors, whether a building has cooling facilities, and whether a building was built after 1986. As a part of the development of a home energy audit methodology, Kim (2014) and Kim and Haberl (2015) used the physical significance of regression coefficients in cooling and heating three-parameter change-point models to calibrate simulation models. Sensitivity analyses were conducted to identify the simulation input parameters that were influential for each of three coefficients, namely the baseload, the change-point temperature, and the cooling/heating slope, and these parameters were varied to allow for the adjusting of the coefficients to match the actual energy use.

The physical interpretation of regression parameter estimates allows one to quantify building parameters that are difficult to measure directly. Many researchers have attempted to use the linear regression models for whole-building energy models to

estimate heat-loss coefficients or overall UA values. However, the accuracy of the estimations varies. As CIBSE (2006) emphasizes, energy regression modeling is not a precise science, and interpretation must be treated carefully; nevertheless, it can draw attention to trends and anomalies, which serve as a starting point for physical investigations that can explain them when there is no detailed knowledge of a building.

1.5.4 Energy-Balance-Based Variables

The key to the practicality of the E_{BL} is that the cooling and heating energy use are canceled out, and the value can represent the net heat load. The energy analysis technique that utilizes the differential of cooling and heating energy use is called an energy-balance-based technique here. The use of this technique is found in some studies.

The first appearance of such a variable in the reviewed literature is in the work of Childs, Courville, and Bales (1983). They used a conceptual value, Q_{net} , which is the total heating energy supplied Q_H subtracted by the total cooling energy removed Q_C . The Q_{net} was used to describe the net load in a building during a daily or longer cycles in the discussion of the influence of thermal mass on building energy consumption. The quantification of a variable similar to Q_{net} is found in Reddy et al.'s research (Reddy et al. 1994; T.A. Reddy et al. 1998), which proposed an index—called the energy delivery efficiency (EDE) index—to assess the efficiency of HVAC air-side systems for commercial buildings as the ratio of the building load (Q_B) to Q_{total} . The Q_B , representing the net cooling load, is defined as the daily cooling energy consumption (Q_C) subtracted by the daily Q_H , whereas Q_{total} is defined as the sum of Q_C and Q_H . The EDE varies from

0 to 1, and the HVAC systems in the building have less mixing of cooling and heating as the EDE value grows larger, which is considered to be more efficient. Reddy (1994) fit a four-parameter change-point (4P-CP) model to the Q_B variable, and the author estimated heat-loss coefficients. The estimated parameters and monitored lights and receptacles data were used to construct an ideal EDE index, which was compared to the actual EDE index to evaluate the performance of HVAC air-side systems.

Deng (1997) and Reddy et al. (Reddy et al. 1999) used the Q_B variable (Reddy et al. 1994; T.A. Reddy et al. 1998) to identify building parameters, including the overall UA value, ventilation rate, indoor temperature, and other non-physical parameters that represent sensible and latent internal loads. The mathematical model of the Q_B variable was formulated, and the parameters were identified using regression estimates. The evaluation of the method used daily Q_B values calculated from simulated hourly system cooling and heating coil loads for one year. A special regression procedure, called a multi-step identification, was proposed, in which correcting the explanatory variable using the estimate from the previous step demonstrate improved bias errors, compared to a regular single-step regression, possibly due to reducing collinearity between the explanatory variables. The parameter identification process was found to be accurate when daily data over an entire year are used. The multi-step identification scheme proved to be accurate; the bias from the true value was less than 10%.

White and Reichmuth (1996) developed a generalized physical model for monthly energy consumption based on the assumption that (1) the difference between the monthly heating and cooling loads and average monthly temperature has a linear

relationship, and the slope is equal to the conduction UA and infiltration heat-loss coefficient of the building, and (2) adding the solar gain and internal gain to this difference of heating and cooling loads results in a line that intersects the average temperature axis at the mean interior set-point temperature of the building (70°F). This physical model was used in an automated tool to assess commercial buildings' energy performance in Reichmuth and Turner (2010) and Reichmuth and Egnor (2013). The physical model was fitted to monthly utility bills, and it was solved by the steepest decent algorithm. The solution set of key physical parameters includes the internal gain, cooling efficiency, service water heating, the heat intercept, the cool intercept, and some other assumed parameters, and they are used for benchmarking and defining various performance indicators to flag poor or unusual performance in individual buildings. This modeling and analysis method was applied to over 3,000 commercial and residential buildings in multiple projects, confirming its broad applicability (Reichmuth and Egnor 2013).

1.5.5 Previous Work on Data Screening Using the E_{BL}

The mathematical description of E_{BL} was derived, and its dependency on the influential building and HVAC control parameters was studied by Shao (2006). Four types of air-side systems, namely a single-duct, constant-air-volume with terminal reheat (CVRH) system; a dual-duct, constant-air-volume (DDCV) system; a single-duct,

variable-air-volume (SDVAV) system; and a dual-duct, variable-air-volume (DDVAV) system, were analyzed using simplified cooling and heating coil load simulations with bin weather data, confirming that E_{BL} does not depend on the air-side system type under the same load conditions. The graphical representation of the simulated E_{BL} versus T_{oa} for varying parameters revealed how the influential parameters, such as outside airflow rate and cooling coil discharge air temperature, can change the slope and intercept.

The data-screening method using E_{BL} , which Shao (2006) proposed, compares the simulated E_{BL} with the measured E_{BL} . To develop the screening threshold, the simulation uncertainty due to the variation in input parameters was estimated using a multivariate analysis for the fractional factorial design simulations. The annual root mean squared error (RMSE), comparing the simulated data with the measured data, was used as a response variable for the multivariate analysis. The most influential factors were determined to be ventilation rate, cold deck temperature, and room temperature; therefore, the variances for these three factors were propagated. Some other uncertainty sources that were not included in the simulation model, such as solar radiation, wind, and occupancy load effects were then added. Based on the estimated uncertainty, the 95% confidence intervals were constructed. If the difference between the simulated and measured data is larger than the interval, then the data are considered to be faulty. The case studies demonstrate that the method is able to identify the outliers, possible scale problems in the measurement, and some types of operational changes. The difficulty in the application of this method lies in obtaining accurate information for the simulation and in estimating the simulation uncertainty. An on-site audit and measurements are

necessary to obtain confident simulation results. The uncertainty of simulation involves numerous factors, including measurement errors, inadequate models, and the use of table values; there is no simple method to determine the total simulation uncertainty.

The *EBL* graphical tool, developed by Baltazar et al. (2007; 2012), tries to achieve a similar level of anomaly detection capability as that of Shao (2006) without using simulations to analyze a large number of buildings every month. The measured data are compared with historical data to find outliers. An example of the energy balance plot used in Baltazar et al. is presented in Figure 1-1. The plot consists of time series of *EBL* and three energy meters—electricity, chilled water, and heating hot water—and scatter plots of the measurements as a function of the T_{oa} . Analysts find outliers or abnormal patterns in the *EBL* vs. T_{oa} plot in the bottom right, and they narrow down the time and meter using the rest of the plots. The energy analysis team in the Energy Systems Laboratory used this graphical tool to assess the quality of energy consumption data collected from over 150 buildings on the Texas A&M University campus, demonstrating the value of the *EBL* vs. T_{oa} plot in the remote detection of anomalies in energy data without detailed information on building use and operation. Two of the more common problems that were detected are related to the scale of the recorded data and to trouble due to the apparent malfunction of the sensors. The errors related to scale factors are typically due to errors in the database processing or in the factors that are set in the data loggers. On the other hand, sensor malfunctioning is a more difficult error to detect (Baltazar et al. 2012).

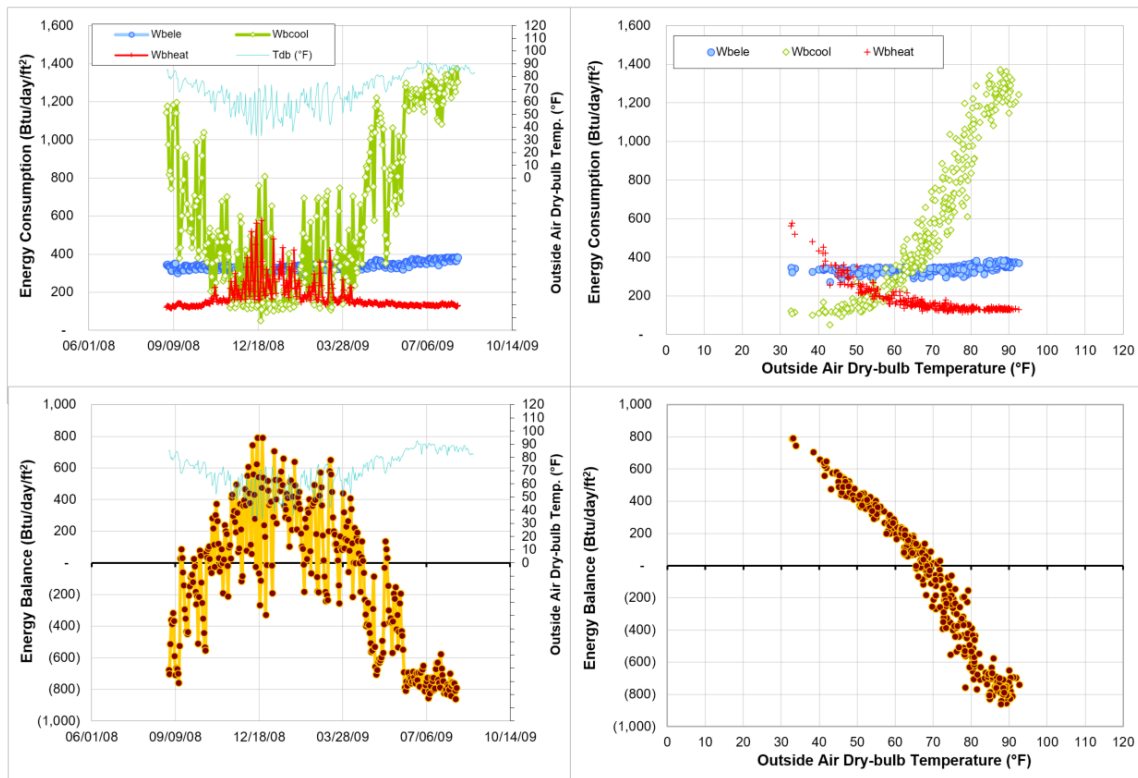


Figure 1-1. Typical energy balance plot for data screening.

To isolate the E_{BL} data points that suggest meter problems, including sensor malfunctions, Masuda et al. (2008) proposed control limits for the E_{BL} vs. T_{oa} plot. Their method applied the four-parameters change-point (4P-CP) model (Kissock et al. 2004) to estimate an E_{BL} regression model as a function of T_{oa} and its statistical uncertainties. To overcome unequal variance or heteroscedasticity due to the presence of outside air latent load, they estimated variable control limits based on the variance of residuals for each temperature bin. As illustrated in Figure 1-2, the resulting control limits for the E_{BL} vs. T_{oa} plots have increasing widths as the T_{oa} increases.

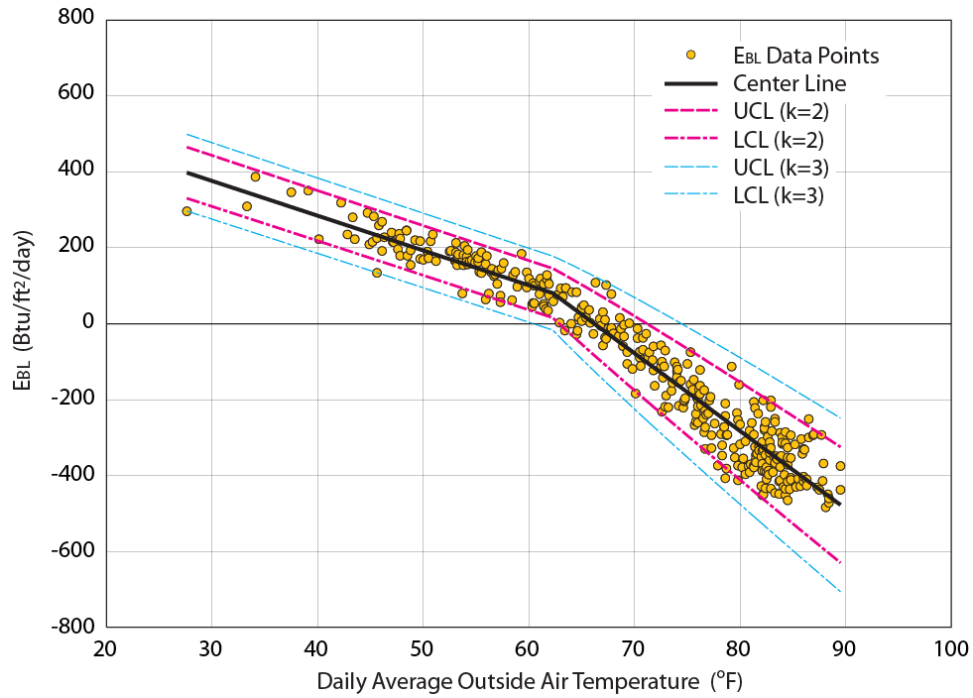


Figure 1-2. Control limits of E_{BL} (Masuda et al. 2008).

The use of outside air enthalpy h_{oa} as an explanatory variable instead of T_{oa} can decrease the degree of heteroscedasticity and the root mean square error (RMSE), especially for high ventilation buildings, allowing one to use constant control limits for annual E_{BL} data (Masuda et al. 2009). Masuda et al. (2009) compared E_{BL} control limits using T_{oa} and h_{oa} for three buildings: an office, a lab, and a dormitory. By using h_{oa} , the RMSEs decreased by 43% and 33% for the office and lab buildings respectively, but increased by 14.6% for the dormitory building, compared to the T_{oa} models. If reliable outdoor humidity data are available in addition to dry-bulb temperature, the use of E_{BL} vs. h_{oa} plots is a viable choice for screening unusual energy data for commercial buildings.

Ji et al. (2008) investigated the sensitivity of building factors on the slope and intercept of the E_{BL} vs. T_{oa} and E_{BL} vs. h_{oa} simple regression models. A multivariate analysis was conducted using spreadsheet energy simulations with the full factorial design with five factors—outside air ratio, cooling coil leaving air temperature, zone temperature, overall UA value, and occupancy density—based on the results of an uncertainty study by Shao (2006). The result demonstrated that the outside air ratio and overall UA value have by far the strongest positive effects on the intercept, followed by the cold deck temperature, in both E_{BL} vs. T_{oa} and E_{BL} vs. h_{oa} models. Similarly, the outside air ratio and overall UA value have strong negative effects on the slope, followed by the zone temperature, in both E_{BL} vs. T_{oa} and E_{BL} vs. h_{oa} models.

1.5.6 Summary of Literature Review

The literature review has presented the vulnerability of thermal energy meters, various energy anomaly detection methods, studies on the physical significance of regression parameter estimates, studies on variables that resemble E_{BL} , and previous works on data screening using E_{BL} .

Thermal energy meters are the least reliable meters among energy meters in buildings because the errors from the flow and temperature measurements are often compounded by calculations. The literature suggests that flow meters can become inaccurate due to various factors, including the inappropriate size selection and installation and off-design conditions, even though the manufacturer's specifications meet accuracy requirements. The errors in thermal energy meters are not easy to detect

because the energy use in a building depends on many factors, including ambient conditions and occupancy. The determination of the influence of these factors in a quick manner is not readily available for data-screening purposes. Therefore, the development of a mechanism to automatically detect errors in thermal energy meters with minimal building information could be highly valuable.

The energy anomaly detection method found in the literature use model-based change-detection that compares model predictions to actual values. For the prediction model, forward modeling based on known physical principles and data-driven modeling using statistical estimation are used. The data-driven models do not require detailed information on building systems, use, and operation; therefore, they are used for applications that involve many buildings. For daily or longer interval data, regression models are widely used because they can be steady, compared to those for the shorter interval data. In this dissertation, the data-driven approach is used because the application is intended for a large number of buildings, and detailed information on each building is not available.

The physical interpretation of energy model regression estimates is an attractive topic for researchers because it allows one to obtain difficult-to-measure values such as overall heat-loss coefficients. Several studies that compared estimated and actual values were reviewed. The estimation accuracy of physical values varies due to the simplification of models and statistical estimation bias. Although the physical interpretation of regression estimates is not an accurate science, several applications successfully used the physical interpretation of parameter estimates for energy audits,

benchmarking, and model calibrations. In this dissertation, the degree of physical significance of the EBL parameter estimates is studied because it can provide useful information to detect physically impossible EBL data-driven models.

The literature review found several examples of EBL -like variables that attempted to represent the building thermal load by subtracting cooling energy use from heating energy use or vice versa. The authors tried to extract information from energy use data as much as possible by using these variables, and they proposed energy performance indicators. This dissertation follows the approach used by Deng (1997) and Reddy et al. (Reddy et al. 1999) for investigating the physical significance of EBL regression parameters.

The previous studies on data-screening that utilized EBL include a forward approach using spreadsheet simulations, a graphical tool for visual anomaly detection, and the development of control limits using the 4P-CP regression models with the outside air temperature or enthalpy as the independent variable. The EBL presents a largely linear pattern as a function of the T_{oa} , and if the building use and operations are roughly constant, the constancy of the pattern can be used to identify abnormal energy data. This technique was used in the monthly verification process for energy use data collected from over 150 buildings on the Texas A&M University campus, where the chilled water and heating hot water are supplied from the central plant, and the while building electricity, chilled water, and heating hot water thermal usage are metered at the whole building level. In this application, the EBL vs. T_{oa} graph has been shown to be a versatile tool for detecting errors in energy use. To detect out-of-pattern EBL data without

relying on a visual analysis, statistical control limits based on the prediction errors of regression models were proposed. This dissertation extends the regression model-based control limits to a continuous and automated monitoring application by using the recursive least squares (RLS) algorithm with a forgetting factor and the Cumulative Sum (CUSUM). As this new approach uses the same linear E_{BL} model as the regression approach, it will automatically update the model whenever new data arrives, and exponentially forget the older data. The new approach addresses the challenge found in the regression approach, namely that the model has to be updated to account for the dynamic use and operation of buildings; for example, renovations, updates to mechanical and lighting systems, and changes in space usage and HVAC controls.

1.6 Significance of the Study

This study is significant because it develops and tests an inexpensive and automated anomaly detection method for whole-facility-level energy use data. The method can be valuable, especially to users who manage a large number of energy meters for commercial buildings in a campus setting.

2. STATISTICAL MODELING OF THE BUILDING ENERGY BALANCE VARIABLE¹

2.1 Introduction

This chapter proposes and investigates several statistical modeling methods for the E_{BL} , for detecting anomalies and biases in the metered energy data using the inverse approach without prior knowledge of the building parameters. The E_{BL} variable is formulated based on the first law of thermodynamics or the energy balance of the building, and it is calculated from separately metered, whole-building daily electricity, cooling, and heating energy consumptions. The use of the E_{BL} variable in addition to the energy data provides an intra-experimental comparison (ASHRAE 2010b) involving measured cooling and heating energy and electricity use based on the energy balance. This is valuable, especially for commercial buildings with simultaneous cooling and heating. Statistical modeling reveals the relationships between E_{BL} and influential factors that are difficult to find by visual examination of the E_{BL} versus T_{oa} plots. And it assists energy analysts in judging data. The E_{BL} model structure is derived using simplified engineering principles, so physical interpretation of the model parameters is possible.

The E_{BL} models developed can be used to detect anomalies and level shifts in the whole-building energy use data in two different cases: (1) newly obtained data are

¹ Reprinted in part with permission from *Energy and Buildings*, Vol 77, Hiroko Masuda and David Claridge, Statistical modeling of the building energy balance variable for screening of metered energy use in large commercial buildings, pp. 292-303. Copyright 2014 Elsevier.

compared to the model prediction, based on the past data; and (2) data in a certain span are filtered retrospectively using the fitted model, based on the same span of data. In either case, the validity of the data is checked in terms of the constancy of the E_{BL} models.

This chapter is organized as follows. Section 2.2 provides a derivation of a functional expression of the E_{BL} variable and the basic structure that will be the foundation of the regression models. Section 2.3 describes the data collected from 56 buildings, and it analyzes the correlations between E_{BL} and the explanatory variables. Section 2.4 proposes four regression models of E_{BL} , and Section 2.5 applies these models to the data and then discusses the results and the applicability of these regression models. Finally, remarks about the proposed methods and further challenges for future studies are presented in Section 2.6.

2.2 Model Structure for the E_{BL}

2.2.1 Model Assumptions

The conditions presented in this section are assumed to simplify the model structure. The purpose of the simplification is to find a basic linear model structure for the daily basis E_{BL} variable that can be applied to the majority of buildings on the Texas A&M University campus in College Station, TX.

2.2.1.1 Indoor air temperature T_{in} is constant

The daily whole-building average of the indoor air temperature is assumed to be constant. Buildings may have different values of T_{in} for occupied and unoccupied hours; however, the daily averages are still roughly constant for most of the campus buildings.

2.2.1.2 Latent load is present when the daily average outside air humidity ratio W_{oa} exceeds $0.01 \text{ kg}_w/\text{kg}_{da}$

This assumes the use of a cooling coil with the leaving air temperature around 12.8°C (55°F) for dehumidification. When the daily average outside air humidity ratio is above $0.01 \text{ kg}_w/\text{kg}_{da}$, the cooling coil latent load tends to be present throughout the day for buildings in College Station, TX. No humidification is provided in any of the campus buildings, so it is not provided in this model.

2.2.1.3 Overall thermal transmittance of the building envelope UA_s is constant

The combined thermal transmittance of the respective areas of gross exterior walls and roof assemblies is assumed to be constant. The heat transfer to the ground is assumed to be negligible, compared to the heat transfer through above-grade envelope assemblies.

2.2.1.4 Daily average of the air exchange rate is constant

The air exchange rate consists of ventilation and infiltration. Commercial buildings have mechanical ventilation to provide the required amount of fresh air for the occupants, and the daily average air exchange rate may be assumed roughly constant. If the building has demand-based ventilation controls, then the variability of E_{BL} between high and low occupancy days may increase. In simplified load calculations, it is common to assume no infiltration in a commercial building during the operation of HVAC systems, assuming that the indoor pressure is maintained at a positive level to prevent infiltration. However, the variation of local pressures caused by various factors, including wind and stack effects in the building, often causes infiltration in real

buildings. Emmerich et al. (2005) reports that infiltration is responsible for 33% of the total heating energy use, but it saves 3.3% of the total cooling energy use in U.S. office buildings, based on dynamic building energy simulation with air flow network modeling. If infiltration is present and turns into cooling/heating loads, it affects the energy balance load and may increase the model errors.

2.2.1.5 Transient effect is negligible

The driving factors of building load such as outside air temperature, solar irradiance, and occupancy schedules, have diurnal variations. These diurnal variations are mostly averaged out in the daily interval data. Some transient effects from these factors remain in the daily data; however, they are assumed to be negligible. The effects from zonal load variations are also assumed to be averaged out in the daily interval energy use data.

2.2.1.6 No economizer and/or heat recovery are present

The impact of the use of an economizer on the whole-building energy use depends on many factors, including outdoor air conditions, AHU types, control strategies, the accuracy of air temperature measurements, and the pressure balance of the building; therefore, it is difficult to develop a whole-building energy model involving economizers without knowing the system's details. College Station, TX, is in a hot and humid climate where economizers are not required in the current building energy code (ICC 2009), and they are not widely used. The use of economizers in the sample buildings is limited and not considered in the models used in this study. Those interested

in the impact of economizers and heat recovery on E_{BL} may wish to examine Shao and Claridge (2006) and Shao (2006).

2.2.1.7 Approximation of solar and occupancy loads

Flouquet (1992) demonstrated that the omission of the solar variable from energy signature models for heating energy use causes serious bias in the parameter estimates. The experiment in (Flouquet 1992) using simulated data demonstrated that the addition of the solar irradiance variable with a solar aperture parameter leads to smaller bias and better estimation of the model parameters in daily, weekly, and monthly frequency data. Sjögren et al. (2009) use a solar variable, which was calculated from a theoretical estimate of the effective solar radiation per square meter, based on the available daily solar irradiance data on the 15th of every month, with 10° of shading and known window areas and orientations on monthly basis energy signature models for heating energy use. Using data from October to March, when the solar radiation is fairly small in Sweden, the estimation of the heat-loss coefficient with the solar irradiance included is slightly smaller for all nine buildings, compared to the estimates without it; however, the impact was fairly insensitive ($\pm 5\%$ for most of the buildings). The indoor temperature estimates were found to be more sensitive to the omission of the solar variable, and estimates increase by 1.9°F on average for all buildings with the solar variable. The multiple linear regression (MLR) study of cooling energy use by Katipamula et al. (1998) found that the addition of global solar irradiance, combined with other weather variables, does not improve the model fit for five large commercial buildings located in central Texas using hourly, daily, monthly, and hour-of-day data.

These studies indicate that (1) the omission of a solar variable may cause the parameter estimates to deviate from the physical values, (2) a solar variable using the window areas and orientations provides reasonable parameter estimates, and (3) the solar variable does not seem to have a strong impact on the prediction errors for cooling energy use. In the present study, window information is not available, and our focus is on prediction rather than parameter estimates. Therefore, we assume a simple linear relationship between the solar load and the solar irradiance, as in Katipamula et al. (1998).

The heat load due to occupants Q_{occ} is expressed as a function of the daily total whole-building electricity use E_{ele} because the diversity of the electricity use has a strong correlation with the diversity of the occupancy level in commercial buildings. For example, Abushakra et al. (2000) deduces that occupancy level in a commercial building as a linear function of the electricity use where direct measurement is not feasible. Q_{sol} and Q_{occ} are modeled as follows:

$$Q_{sol} = a_{sol} + b_{sol}E_{sol} \text{ and } Q_{occ} = a_{occ} + b_{occ}E_{ele} \quad (2.1)$$

where E_{sol} is the daily average global solar irradiance, E_{ele} is the daily total whole building electricity use, and a and b are constants.

2.2.1.8 Model utility

In spite of the limitations imposed by these simplifying assumptions, the E_{BL} model has been shown to be of considerable value in screening metered energy use data from several hundred meters over the last 10 years (Baltazar et al. 2007; 2012). Several campus buildings that are not included in this study demonstrated clustered E_{BL} patterns

due to AHU scheduling. If AHUs are turned off over weekends, then there will be a significant difference in the daily average T_{in} and in the daily average air exchange rate between weekdays and weekends. For such a case, one can group data based on the day types and estimate the model separately for each of the day types (Masuda and Claridge 2012).

2.2.2 Formulation of E_{BL} Model Structure

Each of the terms on the right-hand side of Equation (1.2) can be expressed as presented in this section, based on the simplified load calculation principles and the assumptions discussed above. The air exchange load Q_{air} is the total of the sensible and latent air exchange load, which is as:

$$Q_{air} = m_v c_p (T_{oa} - T_{in}) + m_v h_v W_{oa}^+ \quad (2.2)$$

where m_v is the air exchange rate, c_p is the specific heat, h_v is the specific heat of vaporization, and W_{oa}^+ is the humidity load variable, which is defined as $W_{oa}^+ = (W_{oa} - W_{threshold})^+$. The outside air humidity ratio threshold $W_{threshold}$ for the presence of latent load is set to 0.01 kgw/kgda based on the assumption of Section 2.2.1.2. The superscript ‘+’ indicates that the term is used only if it is positive; otherwise it is 0. The building envelop conduction load is as:

$$Q_{cond} = UA_s (T_{oa} - T_{in}) \quad (2.3)$$

where UA_s is the overall thermal transmittance of the building envelope. By inserting Equations (2.1)–(2.3) into Equation (1.2), letting $f = 1$ and organizing the parameters by the measurable variables, the E_{BL} model structure is derived as:

$$E_{BL} = \beta_0 + \beta_T T_{oa} + \beta_W W_{oa}^+ + \beta_{sol} E_{sol} + \beta_{occ} E_{ele} \quad (2.4)$$

where each β is an unknown parameter. It should be emphasized that the parameters β_T and β_W have physical significance; $\beta_T = -(UA_s + m_v c_p)$ and $\beta_W = -m_v h_v$. The implication of the intercept is $\beta_0 = (UA_s + m_v c_p) T_{in} - a_{sol} - a_{occ}$. The model structure in Equation (2.4) is used as the basis of the regression models in Section 2.4.

2.3 The Data

2.3.1 Description of the Data

The energy use data from 56 sample buildings on the Texas A&M University College Station campus were used to study E_{BL} regression models. The buildings include a variety of functions, for example, offices, classrooms, laboratories, and dormitories, among others, and the gross floor areas range from 323 m² (3477 ft²) to 32,116 m² (345,694 ft²), as listed in Table 1. All of the buildings have separately metered hourly data for electricity, chilled water, and heating hot water energy consumption. The daily consumption data were summed from the hourly data, and the E_{BL} values were calculated from the daily energy consumption using Equation (1.2). The data used are for the 2011 calendar year (January 1, 2011 through December 31, 2011). Prior to the model estimations, some extreme data values were removed by visual inspection using time series plots and E_{BL} versus T_{oa} plots. The removed data are obvious outliers, which are outside the normal groups by approximately 5σ or more, and these data errors are typically caused by sensor failures. Outlier detection for such extreme values can be done statistically using robust regressions or pattern recognition techniques; however, these were not implemented in this study. After this process, 51 of the 56 buildings have

at least 75% of annual daily observations. The gross floor area, the number of daily E_{BL} data points used for modeling, and the main functions of the sample buildings are presented in Table 2-1. The hourly interval observations of dry-bulb and wet-bulb temperatures were obtained from the Quality Controlled Local Climatological Data (QCLCD) at the local weather station (NOAA 2013), and the daily average of the T_{oa} and W_{oa} were calculated from the hourly observations. The 15-minute interval global solar irradiance data were obtained from the solar test bench located on the Texas A&M University campus (Baltazar et al. 2011), from which the daily average data were calculated.

Table 2-1. Gross floor area, number of daily E_{BL} data points, and functions of the sample buildings.

Building number	Area (m ²)	Number of data	Building functions
1	6,251	304	Dormitory
2	6,251	295	Dormitory
3	6,251	316	Dormitory
4	7,689	319	Dormitory
5	11,610	323	Sports
6	14,296	284	Sports
7	10,245	337	Office
8	1,799	297	Office
9	19,045	242	Laboratory, Office, Classroom
10	10,834	330	Office, Classroom

Table 2-1 Continued.

Building number	Area (m ²)	Number of data	Building functions
11	2,968	331	Dormitory
12	3,427	329	Dormitory
13	3,793	358	Dormitory
14	3,793	356	Dormitory
15	4,193	352	Dormitory
16	4,484	332	Office
17	14,339	256	Dormitory
18	3,619	358	Dormitory
19	12,156	358	Office
20	8,623	358	Office
21	19,444	331	Theater, Office
22	8,922	348	Laboratory, Office, Classroom
23	7,670	328	Office
24	2,759	358	Office
25	3,706	299	Office, Classroom
26	3,722	358	Office
27	6,494	357	Office
28	3,687	358	Office, Classroom
29	1,772	343	Office, Classroom
30	11,230	323	Office, Classroom
31	6,444	331	Laboratory, Office, Classroom
32	30,138	334	Office, Classroom

Table 2-1 Continued.

Building number	Area (m ²)	Number of data	Building functions
33	5,882	357	Health Center
34	9,750	339	Classroom, Laboratory, Office
35	23,965	358	Office, Classroom
36	5,774	353	Dormitory
37	6,472	334	Dormitory
38	6,472	261	Dormitory
39	4,364	355	Laboratory, Office
40	9,610	333	Animal Hospital, Office
41	1,600	297	Laboratory, Office
42	13,087	167	Animal Hospital, Office
43	323	355	Office
44	14,770	342	Laboratory, Office, Classroom
45	15,780	358	Laboratory, Office
46	11,023	352	Laboratory, Office, Classroom
47	2,570	358	Laboratory, Office
48	6,329	348	Library
49	3,464	170	Laboratory, Office
50	21,882	296	Laboratory, Office
51	32,116	327	Sports
52	11,304	341	Museum, Office, Archive
53	12,386	338	Office, Classroom
54	6,197	340	Office

Table 2-1 Continued.

Building number	Area (m ²)	Number of data	Building functions
55	2,019	357	Laboratory, Office
56	9,662	350	Laboratory, Office

2.3.2 Correlation Analysis of the Data

The correlations between the E_{BL} variable and each of the variables— T_{oa} , W_{oa}^+ , E_{sol} , and E_{ele} —on the right-hand side of Equation (2.4) were examined for all the sample data using the semi-partial correlation coefficient sr which defines the contribution of each explanatory variable to the multiple correlation. For example, the sr between E_{BL} and T_{oa} is the correlation between E_{BL} and T_{oa} from which the effects of other variables have been removed. The coefficient is normalized to the range from -1 to 1 (positive and negative signs mean positive and negative correlations respectively); therefore, it is useful to compare variables with different units and for different buildings. The derivation of the coefficient is explained in statistical textbooks such as Cohen and Cohen (2003). The coefficients for the data from 56 buildings were estimated using the R package ‘ppcor’ (Kim 2012), and the results are presented in Figure 2-1. The results demonstrate that the T_{oa} variable has the largest effect for all the buildings, followed by the W_{oa}^+ and E_{sol} variables. The buildings with high sr for the W_{oa}^+ variable (buildings 9, 55, 31, 41, and 45) are laboratory buildings that are known to have large ventilation rates, while the buildings with low sr for the W_{oa}^+ variable (buildings 3, 29, 15, 38, and

1) are dormitories, which have relatively low ventilation rates among campus buildings. This matches the functional expression derived in Section 2.2; the parameter of W_{oa}^+ is a function of the air exchange rate m_v of the building. Even though the signs of the parameters in Equation (2.4) and the sr 's are supposed to be negative, the E_{ele} variable has a positive sr for 22 buildings. One plausible reason for these contradictory results is that the electricity consumption level may be related to some factor(s) other than occupancy level. For instance, if the indoor air temperature floats above the set point (for cooling) or below the set point (for heating) on less occupied days, the E_{BL} may exhibit a positive relation to the electricity consumption level. Since the contribution of occupancy load to the total cooling load in a large commercial building is small, these temperature variations may obscure the correlation between the Q_{occ} and E_{ele} . Similar phenomena may occur when the building has demand-based ventilation controls. The physical meaning of the effect of the E_{ele} variable is not clear and may differ for each building; therefore, this variable is not included in the regression models in this study.

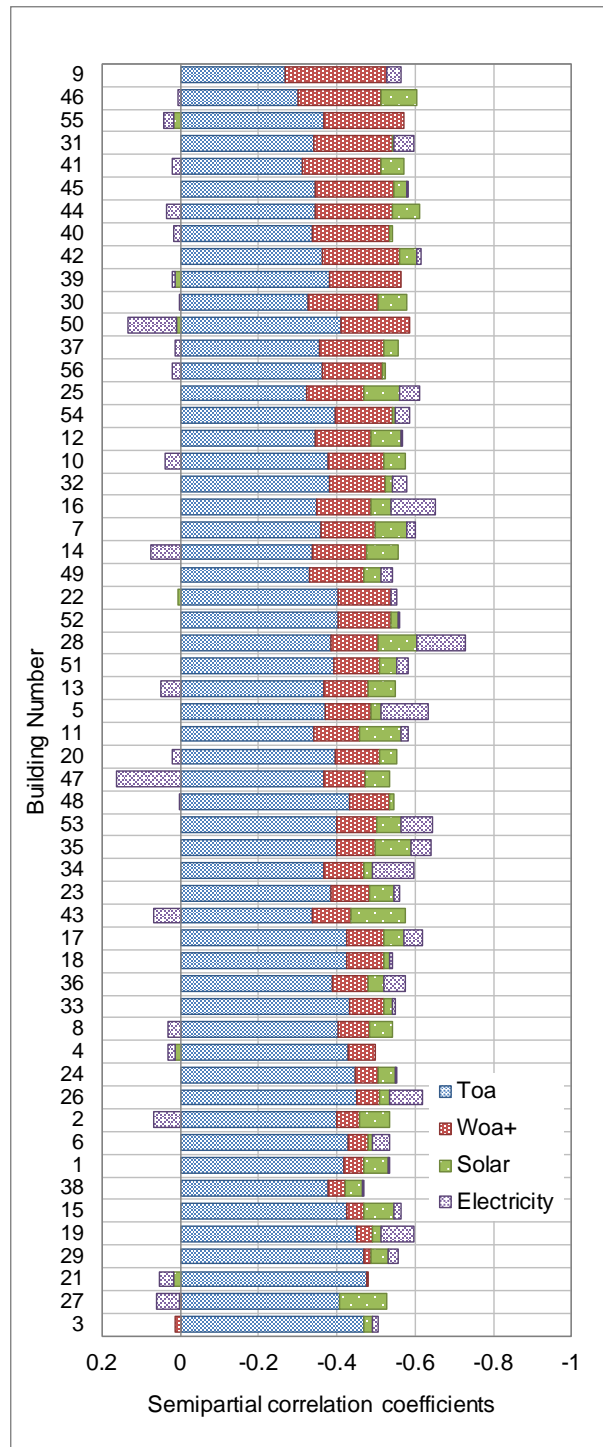


Figure 2-1. Semipartial correlation coefficients sr between the response variable E_{BL} and the explanatory variables Toa , W_{oa}^+ , E_{sol} , and E_{ele} for the sample buildings. The buildings are ordered from the smallest to largest values for W_{oa}^+ .

2.4 Regression Models

2.4.1 The Four-Parameter Change-Point (4P-CP) Model

Schrock and Claridge (1989) and Ruch and Claridge (1992) proposed the four-parameter change-point (4P-CP) model to develop temperature-dependent baselines for large commercial buildings. Kissock et al. (1998) and Reddy et al. (1999) described the physical basis of change-point models, including the 4P-CP model. These modeling procedures are incorporated into IPMVP (DOE 1997; EVO 2012), and they are widely used for baseline modeling to measure the energy savings from retrofit projects. The 4P-CP model consists of two lines with one break point, which is called the change-point temperature T_{cp} , and it can be written as:

$$E_{BL} = \beta_{cp} - \beta_1(T_{cp} - T_{oa})^+ + \beta_2(T_{oa} - T_{cp})^+ + \varepsilon \quad (2.5)$$

where T_{cp} is the E_{BL} at $T_{oa} = T_{cp}$, β_1 and β_2 denote the left and the right slopes, and ε is an error term. The superscript ‘+’ means that the value inside the bracket will be treated as 0 when it is negative. The estimation method used in this study is based on the one in ASHRAE IMT (2004). In the estimation procedure, the data are divided at a T_{cp} , and the slope for each of the groups is separately estimated by ordinary least squares (OLS). This procedure is repeated by incrementing the T_{cp} . The best fit model is the model for the T_{cp} value that yields the minimum root mean squared error (RMSE). The search for the best fit is done in two stages; first, using a rough temperature grid, and second, using a finer grid inside the regime found in the first stage. The ASHRAE IMT includes the Fortran 90 source code and the test report that demonstrates the stability and accuracy of the algorithm (2004).

2.4.1.1 *Non-constant error variance in 4P-CP models*

When the 4P-CP models are fitted to E_{BL} data, the model residuals are prone to increase with the outside air temperature, especially for buildings with high ventilation rates, due to the latent load. With a non-constant error variance, even though the least-squares estimators will still be unbiased, they will no longer have the minimum-variance property. However, the impact of non-constant error variance on the efficiency of the OLS is not always serious. Using a simple case with the error variance proportional to the explanatory variable, Fox (2008) demonstrated that the precision of the OLS estimator stabilizes quickly as the sample grows, and the author concluded with a rough rule that non-constant error variance seriously degrades the least-squares estimator only when the ratio of the largest to the smallest variance is approximately 10 or more. The variation of the error variance in E_{BL} 4P-CP models is typically less than the level described in this rule, and the impact on the model parameter estimates should not be significant. The problem regarding data screening is that the prediction intervals that are estimated based on the constant error assumption provide intervals that are either too wide or too narrow for screening, depending on the level of T_{oa} . The use of outside air enthalpy h_{oa} instead of T_{oa} decreases the degree of unequal variance and provides better prediction in the region where latent load is present (2009), although it tends to increase errors in the region without outside air latent load. This method is preferred when outside air humidity measurements are available but matrix computations are avoided. Another method to handle unequal error variance is to adjust the widths of prediction intervals. Masuda et al. (2008) proposed a method to construct prediction intervals for

the 4P-CP E_{BL} models as a function of T_{oa} based on the local variance estimated for each temperature bin. This approach can be used when humidity measurements are not available.

2.4.2 Multiple Linear Regression (MLR) Models

If humidity measurements are available in addition to outside dry bulb air temperature, one can estimate the multiple linear regression (MLR) model with T_{oa} and W_{oa}^+ :

$$E_{BL} = \beta_0 + \beta_T T_{oa} + \beta_W W_{oa}^+ + \varepsilon \quad (2.6)$$

where β_0 is the intercept, β_T and β_W are the parameters of T_{oa} and W_{oa}^+ respectively, and ε is an error term. Solar irradiance measurements are not available for many locations; however, according to the correlation analysis, inclusion of the global solar irradiance variable E_{sol} (when available) is likely to improve the model fit for many buildings. With this variable, the MLR model is:

$$E_{BL} = \beta_0 + \beta_T T_{oa} + \beta_W W_{oa}^+ + \beta_{sol} E_{sol} + \varepsilon \quad (2.7)$$

where β_{sol} is the parameter of E_{sol} , and ε is an error term. The parameters in the MLR model are estimated by OLS.

2.4.2.1 Multicollinearity in MLR models

The existence of one or more strongly correlated explanatory variables results in large variances and covariances, and it also tends to produce least-squares estimates that are too large in absolute value (Montgomery et al. 2012). This problem is called multicollinearity. The variance indication factor (VIF) is a widely used multicollinearity diagnostic, and it is evaluated as $VIF = (1 - R_i^2)^{-1}$, where R_i^2 is the coefficient of

determination between the i^{th} explanatory variable and all of the other explanatory variables in the regression equation. The VIF is close to unity when multicollinearity does not exist, and it becomes large as the multicollinearity level increases. There is no objective judgment for the value of the VIF at which multicollinearity becomes a serious problem; however, several investigators, including Haan (2002) and Montgomery (2012), refer to the convention that a VIF exceeding 5 or 10 is considered to be an indication of a serious level of multicollinearity that affects the parameter estimates. The VIFs of the E_{BL} MLR models for sample buildings are in the range of 2.4 to 3.1 for the two-variable models (T_{oa} and W_{oa}^+), and they are in the range of 1.6 to 4.8 for the three-variable models (T_{oa} and W_{oa}^+ , and E_{sol}), so the problem of multicollinearity should not be serious. Additionally, the fitted model often produces satisfactory predictions despite the presence of strong multicollinearity, and this level of multicollinearity is not an issue for our purpose of having better E_{BL} predictions for data screening.

2.4.3 *Multiple Linear Models Incorporating an Autocorrelated Error Structure*

Time series data often exhibit serially correlated errors, and such errors are said to be autocorrelated. With autocorrelated errors, the OLS estimates are still unbiased; however, they are no longer minimum-variance estimates, and the standard errors (SEs) are wrong (Montgomery et al. 2012). Positive autocorrelations are often found in errors from linear regression models of daily interval building energy data, as discussed by Ruch et al. (1999) and Reddy et al. (1998). The daily basis E_{BL} , which is evaluated from energy use, also displays autocorrelated residuals when the models are estimated by OLS. Autocorrelated residuals are caused by model misspecifications (Ruch, J. K.

Kissock, et al. 1999). For example, the variations of cooling and heating loads from building schedules and different control settings for weekday/weekend and the thermal inertia effect on the building load are not included in the simplified linear models, and these can cause autocorrelated residuals. The systematic bias of a model caused by the change in the parameters during the modeling period may generate residuals with a periodic pattern inherited from the explanatory variables, which can also be a source of autocorrelated residuals. The application of linear models with autocorrelated error structures for daily building energy use data are found in Ruch et al. (1999) and Liu et al. (2011). Both use OLS estimates for the main model structure, and they fit autoregressive or more advanced time series models to the residuals. Ruch et al. (1999) fitted a first order autoregressive model to the residuals of the baseline model to improve the savings estimation for retrofits, based on the assumption that the baseline and the post-retrofit periods have the same autoregressive parameter. Liu et al. (2011) fitted seasonal autoregressive integrated moving average (ARIMA) models, and they found the best-fitting combination of ARIMA parameters by the Bayesian information criterion (BIC) to have better model fitting for anomaly detection and forecasting. To observe the effectiveness of including autocorrelated error structures in MLR models for the energy balance load, the MLR model that incorporates errors with an autoregressive correlation structure of order 1 (AR[1]) is specified as:

$$\begin{aligned}
 E_{BL,t} &= \beta_0 + \beta_T T_{oa,T} + \beta_W W_{oa,t}^+ + \varepsilon_t \\
 \varepsilon_t &= \rho \varepsilon_{t-1} + v_t
 \end{aligned} \tag{2.8}$$

where ε_t is the error term at time period t , and ρ ($|\rho| < 1$) is the autocorrelation parameter. ε_t depends on the error in the previous time period ε_{t-1} and on a white noise process v_t . The AR(1) structure is assumed for the error term ε_t because the residuals from the OLS estimations of the EBL models in Equations (2.6) and (2.7) for the sample buildings exhibit strong autocorrelations in lag 1 (one day). The Durbin-Watson test for the AR(1) structure with the null hypothesis— $H_0: \rho = 0$ (no autocorrelation present)—is strongly rejected ($p < 0.0001$) for the majority of the sample buildings. The parameters in Equation (2.8) are estimated using the maximum likelihood (ML) method.

2.4.4 Evaluation of the Fitted Models

The model fitting is examined by comparing the actual and fitted values. The RMSE, a measure of the spread of the data around the estimated model, is estimated for this purpose. The RMSE is estimated as:

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - p}}, \quad i = 1, 2, \dots, n \quad (2.9)$$

where y_i is the i th observation (in our case, EBL), \hat{y}_i is the fitted value, n is the number of observations, and p is the number of parameters in the model. The approximate uncertainty intervals are estimated as (fitted value) $\pm k \cdot \text{RMSE}$, where k is a multiplicative factor to adjust the desired uncertainty level. Since daily data for a year are used in this study, the number of observations is large, and the residuals tend to follow normal distributions. For $k = 2$, 95.4% of the data will be within the intervals, and for $k = 3$, 99.7% of the data will be within the intervals. To compare the degree of overall model fit between buildings, the coefficient of variation (CV) of the residuals is

estimated. The CV represents the ratio of the deviation to the sample mean. The CV values for E_{BL} models that are estimated as Equation (2.9) become large or negative because the annual mean of E_{BL} in a hot and humid climate is both close to 0 and often negative. Therefore, we use the magnitude of the E_{BL} variation to normalize the deviation instead of the mean. The CV for the E_{BL} model CV-RMSE is calculated as:

$$\text{CV-RMSE} = \frac{\text{RMSE}}{[(E_{BL,Max} - E_{BL,Min})/2]} \quad (2.10)$$

where RMSE is the root mean square error, and $E_{BL,Max}$ and $E_{BL,Min}$ are the maximum and minimum E_{BL} values in the dataset respectively. The better the model fit, the smaller the CV-RMSE becomes. Building energy use data sometimes have shifted levels due to sensor bias and errors in calculation constants, among other reasons. These bias errors cannot be detected as a form of outliers if all the past data are biased as well. We learned that the outside air temperature in which $E_{BL} = 0$ is a valuable reference to assess the validity of E_{BL} levels. Based on the E_{BL} versus T_{oa} plots for over 100 buildings on the Texas A&M University campus during the past several years, this temperature varies between buildings, but it generally distributes around 70°F (21.1°C), where the daily average indoor temperatures of the campus buildings are presumed to be in the range of 73–76°F (22.8–24.4°C). When this temperature is out of the range of 60–80°F (15.6–26.7°C), we often find biased readings such as errors in calculation constants. This temperature is denoted as the indoor reference temperature T_{ref} because the major factor determining this temperature is the indoor air temperature. To estimate the T_{ref} for

multivariate models, we define T_{ref} as the negative of the intercept β_0 , divided by the temperature slope β_T . The physical significance can be derived as:

$$\begin{aligned} T_{ref} &= -\frac{\beta_0}{\beta_T} \\ &= T_{in} - \frac{Q_{occ} + Q_{sol}}{m_v c_p + UA_s} \end{aligned} \quad (2.11)$$

Note that Q_{sol} will be mitigated if the solar irradiance variable E_{sol} is included in the model. For the 4P-CP models, the T_{ref} is estimated using the left slope if $T_{cp} < 0$, otherwise the right slope is used. The T_{ref} resembles the balance-point temperature, which is used as the base temperature of the VBDD method to calculate building heating and cooling loads. The major difference is that T_{ref} does not include the internal load from electricity, but it does include air exchange load. This difference makes T_{ref} more stable between buildings with different load characteristics, compared to the balance point temperature. From the expression in Equation (2.11), one can deduce some attributes of T_{ref} . The value of T_{ref} approaches T_{in} as the air exchange rate of the building becomes large, and buildings with a high Q_{occ} or Q_{sol} may have lower T_{ref} .

2.5 Results and Discussion

The four models described in Sections 2.4.1–2.4.4, namely (1) the 4P-CP model with T_{oa} ; (2) the MLR model with T_{oa} and W_{oa}^+ ; (3) the MLR model with T_{oa} , W_{oa}^+ , and E_{sol} ; and (4) the multivariate linear model with T_{oa} and W_{oa}^+ , incorporating the AR(1) error structure, were all estimated using the data from the 56 sample buildings, and the T_{ref} values were calculated from these parameter estimates. Estimations for the 4P-CP models were performed using an in-house program based on the algorithm in ASHRAE

IMT (Kissock et al. 2004), and estimations for multi-variable models were performed using R (R Core Team 2013). In the following sections, these models are designated as 4P-CP(T), MLR(T,W), MLR(T,W,S), and AR1(T,W) respectively.

2.5.1 Comparison of the Model Fits

The mean values of CV-RMSE, calculated according to Equation (2.10), are 10.4% (SE = 0.34%), 10.0% (SE = 0.41%), 9.7% (SE = 0.39%), and 6.9% (SE = 0.22%) for the 4P-CP(T), MLR(T,W), MLR(T,W,S), and AR1(T,W) model groups respectively, where SE is the standard error of the mean. Figure 2-2 illustrates the distribution of the CV-RMSE values for each model group as a box-whisker-mean plot. The degree of overall model fit is approximately the same for the first three models, but significantly better in the AR1(T,W) models. Inclusion of the E_{sol} variable does not have a large impact on the model fit for most of the buildings.

While the 4P-CP(T) and MLR(T,W) models have a similar level of CV-RMSE on average, the results for individual buildings are quite different. Figure 2-3 presents the CV-RMSE values for the 10 buildings with the largest srs for the W_{oa}^+ (i.e. high ventilation) in Figure 2-1, and Figure 2-4 illustrates the results for the 10 buildings with the lowest srs for W_{oa}^+ (i.e. low ventilation). These plots indicate that the MLR(T,W) models provide significant improvement, compared to the 4P-CP(T) models, for the high ventilation buildings, but not for the low ventilation buildings. This is due to the decrease in unequal residual variances by including the humidity variable. The plot of residuals against T_{oa} for building 31 is presented in Figure 2-5 as an example. This plot also indicates that exclusion of the E_{sol} variable does not cause serious systematic and

skewed errors, unlike the W_{oa}^+ variable. Considering that the solar variable is not readily available and that the variable does not have a large impact on the model fit for this sample of buildings, the use of the T_{oa} and W_{oa}^+ variables may be practically sufficient for multiple linear E_{BL} models.

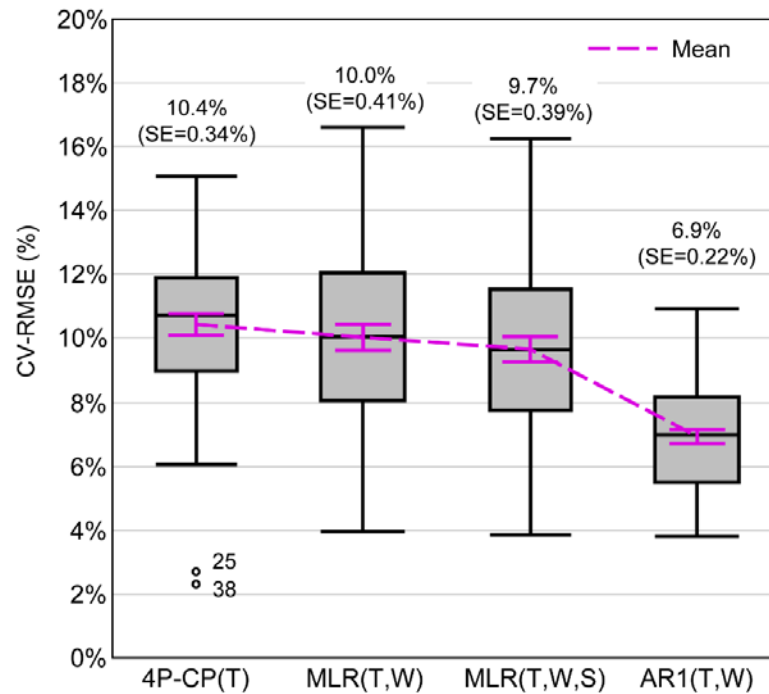


Figure 2-2. The CV-RMSE values of the fitted models using data from the 56 sample buildings. A box represents the interquartile range (IQR) with the median as a horizontal line, and the whiskers extend from the ends of the box to the outermost data point that falls within the 75th quantile +1.5·IQR and the 25th quantile -1.5·IQR. The data outside this range are presented as outliers, and the building numbers are indicated. The broken line connects the means, and the inside error bars are constructed using 1 SE from the mean. The mean and SE values of each model are indicated in the plot.

The 4P-CP(T) models for buildings 25 and 38 have small errors (CV-RMSE = 2.7% and 2.3% respectively), which appear as outliers in Figure 2-2. The parameters for these buildings seem to have changed during the course of the period that

was modeled, and the MLR models, based on a constant model structure, could not accurately represent the data as a whole. Meanwhile, the 4P-CP(T) model does not require the left and right slopes to have the same physical properties, and the model was empirically fitted to have the smallest residuals. Although the 4P-CP(T) has a better fit, the structure change during the modeling period, which might be important for energy analysis, is not revealed. In MLR models, changes in the model structure appear as systematic residual behavior, as illustrated in Figure 2-6, which is a residuals versus T_{oa} plot for building 25. There are many possible causes for this type of systematic residual pattern, including a level shift in the sensor readings used for energy metering, planned or unplanned operation and control changes, and a retrofit of the building. Bias in the metered energy consumption also causes a similar residual pattern; for example, when the chilled water temperature difference between supply and return is out of calibration and is measured to be larger than the actual values, the unbalance of the E_{BL} leads to a skewed pattern in the high temperature region where cooling energy consumption becomes large.

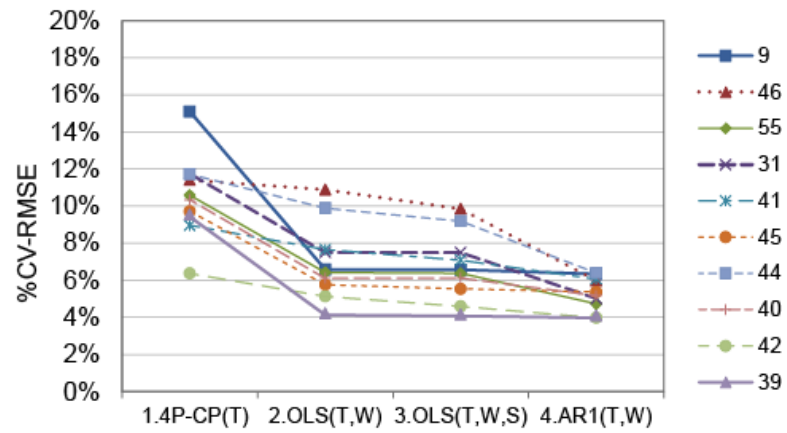


Figure 2-3. The CV-RMSEs of the fitted models for the 10 buildings where the W_{oa}^+ variable had the greatest impact.

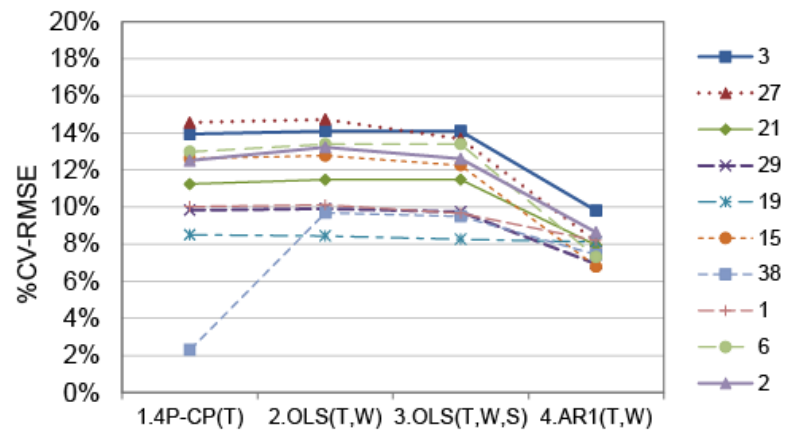


Figure 2-4. The CV-RMSEs of the fitted models for the 10 buildings where the W_{oa}^+ variable had the least impact.

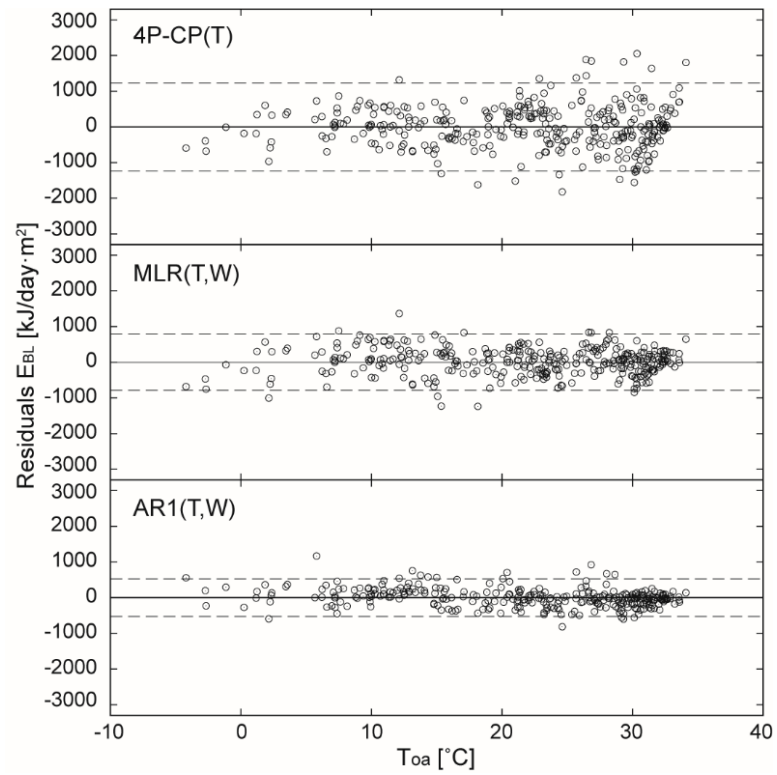


Figure 2-5. E_{BL} residuals versus T_{oa} for 4P-CP(T), MLR(T,W), and AR1(T,W) models for building 31. The dashed lines represent approximate uncertainty intervals at $k = 2$ for the respective models.

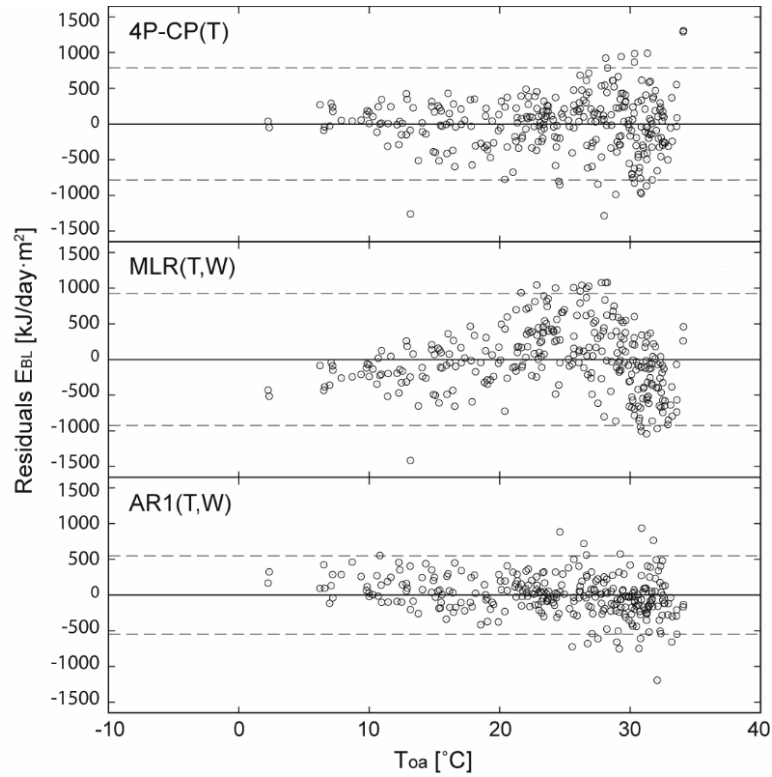


Figure 2-6. The E_{BL} residuals versus T_{oa} for the 4P-CP(T), MLR(T,W), and AR1(T,W) models for building 25. The dashed lines represent approximate uncertainty intervals at $k = 2$ for the respective models.

2.5.2 Inclusion of AR(1) Error Structure in the MLR Model

Inclusion of the AR(1) error structure in the MLR(T,W) model decreased the CV-RMSE from 10% to 6.9% on average. The estimated autocorrelation coefficients for the AR1(T,W) models were in the range of 0.23–0.97. There are two types of improvements in the model fit; one is related to outside air temperature variations, and the other is not. The first type is possibly associated with the lag and damping of the heat load due to the thermal inertia of the building mass. Figure 8 illustrates the time series of the actual E_{BL} and the fitted values by the MLR(T,W) and the AR1(T,W) for building 7. The daily average T_{oa} in College Station has large fluctuations during the winter season,

from November through to February, every year. A similar temperature pattern is observed in Houston, TX, and this appears to be a weather characteristic of this region. These fluctuations are not filtered out by daily averaging; therefore, the frequency is lower than the 48 h^{-1} based on the Nyquist frequency $f_N = 1/(2\Delta t)$ for a sampling interval $\Delta t = 24 \text{ h}$. The OLS models over-fitted the response to this long-period T_{oa} variation for some buildings, including building 7; however, incorporating the AR(1) error structure into the models appeared to better capture this damping effect.

The second type of improvement may be related to operational changes in the building that are not related to T_{oa} variations. Level changes over weekends and during breaks were better fitted by the AR(1) error structure. Other variations, whose reasons are unknown, were also fitted better. For example, building 11 had an increase of E_{BL} for 5 days during the period of July 1721, 2011, caused by a sudden drop in the chilled water flow rate. As depicted in Figure 9, the AR1(T,W) model captured this increase well, although we want to flag such a disruption in energy data. The inclusion of the AR(1) error structure in the MLR model may better fit E_{BL} variations that are unspecified in the model; however, this could be a drawback at the same time because the cause of the level change might be metering problems.

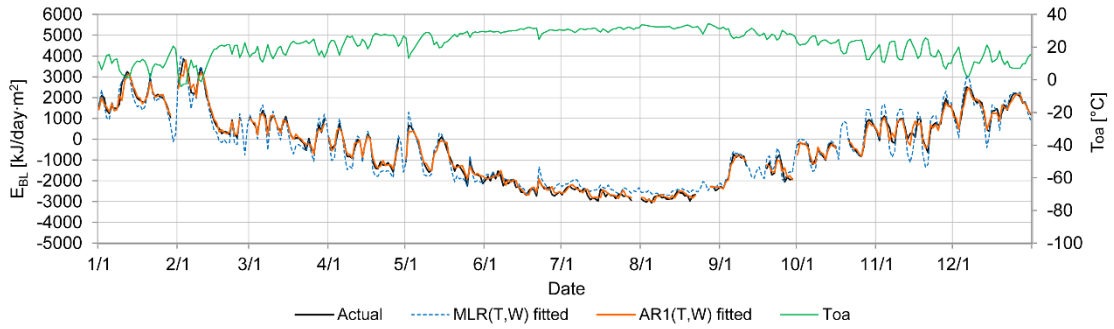


Figure 2-7. Time series plots of the T_{oa} , the actual E_{BL} , the fitted E_{BL} by the MLR(T,W) model (CV-RMSE = 9.7%), and the fitted E_{BL} by the AR1(T,W) model (CV-RMSE = 5.4%) for building 7.

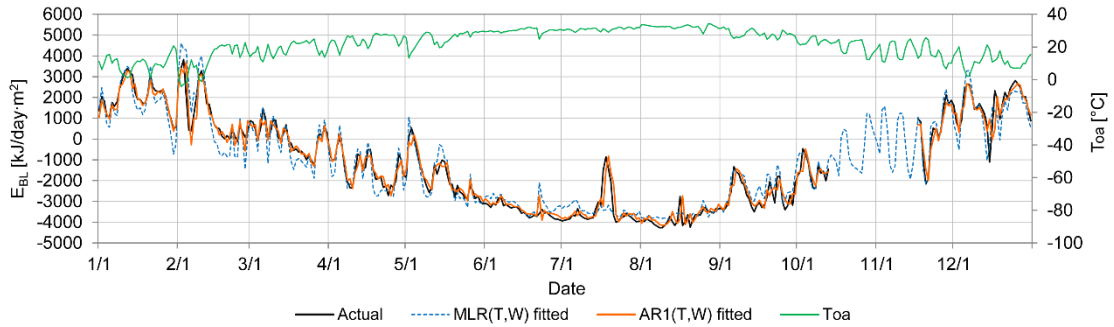


Figure 2-8. Time series plots of the T_{oa} , the actual E_{BL} , the fitted E_{BL} by the MLR(T,W) model (CV-RMSE=13.8%), and the fitted E_{BL} by the AR1(T,W) model (CV-RMSE = 9.3%) for building 11.

2.5.3 Comparison of the Indoor Reference Temperature

The means of the T_{ref} values are 19.1°C (SE = 0.24°C), 19.3°C (SE = 0.27°C), 20.9°C (SE = 0.31°C), and 17.7°C (SE = 0.48°C) for the 4P-CP(T), MLR(T,W), MLR(T,W,S), and AR1(T,W) models respectively, as illustrated in Figure 2-9, where SE is the standard error of the mean. The first two models do not demonstrate a significant difference in the means. The inclusion of E_{sol} variable increased the average T_{ref} by 1.6°C from the MLR(T,W) models, which can be explained from Equation (2.11) with a

smaller Q_{sol} in the numerator. Building 2 is presented as an outlier in the 4P-CP(T), MLR(T,W), and AR1(T,W) model groups, but not in the MLR(T,W,S) model group, since the E_{sol} variable has a strong effect on the E_{BL} in the data for building 2. The AR1(T,W) model demonstrates a lower average T_{ref} , compared to the other groups. This is because the absolute values of the parameter estimates in the AR1(T,W) models are more conservative than the absolute values of the OLS estimates. The use of T_{ref} for checking the validity of the E_{BL} level for the AR1(T,W) models may not be as effective as for other OLS models because the T_{ref} values in the AR1(T,W) model group are spread out more widely than in other model groups.

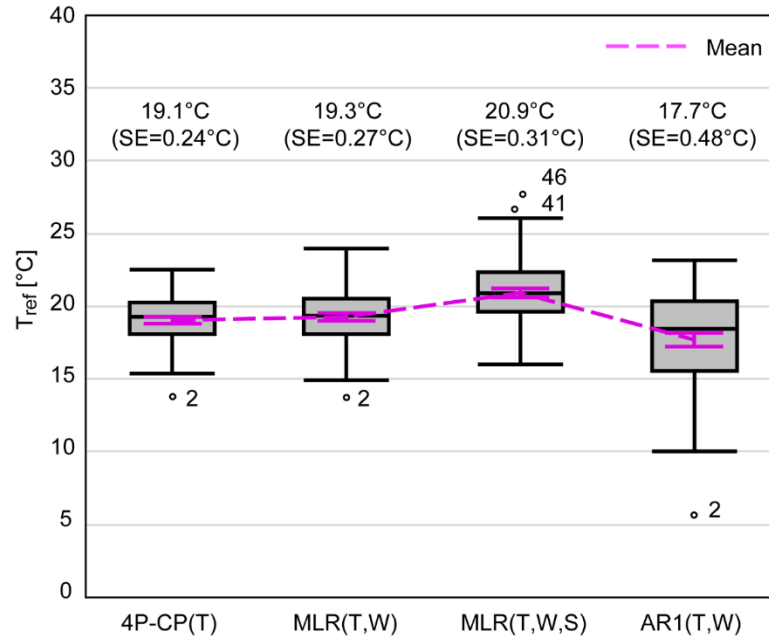


Figure 2-9. The T_{ref} of the fitted models using data from the 56 sample buildings. A box represents the interquartile range (IQR) with the median as a horizontal line, and the whiskers extend from the ends of the box to the outermost data point that falls within the 75th quantile +1.5·IQR and the 25th quantile -1.5·IQR. The data outside this range are presented as outliers, and the building numbers are indicated. The broken line connects the means, and the inside error bars are constructed using 1 SE from the mean. The mean and SE values of T_{ref} for each model are indicated in the plot.

2.5.4 Change Point Temperature

The histogram of the T_{cp} values estimated from the 4P-CP(T) models is plotted with the daily average weather data in Figure 11. The T_{cp} in a 4P-CP E_{BL} model represents the T_{oa} where the latent load starts to appear, which is equivalent to the W_{oa}^+ threshold for E_{BL} multivariate models. The median T_{cp} for the sample data is 19.5°C (66°F), and it separates the outside air condition into approximately the same two groups as the $W_{oa} = 0.01 \text{ kg}_w/\text{kg}_{da}$ line. This indicates that the T_{cp} in the 4P-CP(T) model and $W_{threshold}$ in the MLR models have essentially the same meaning.

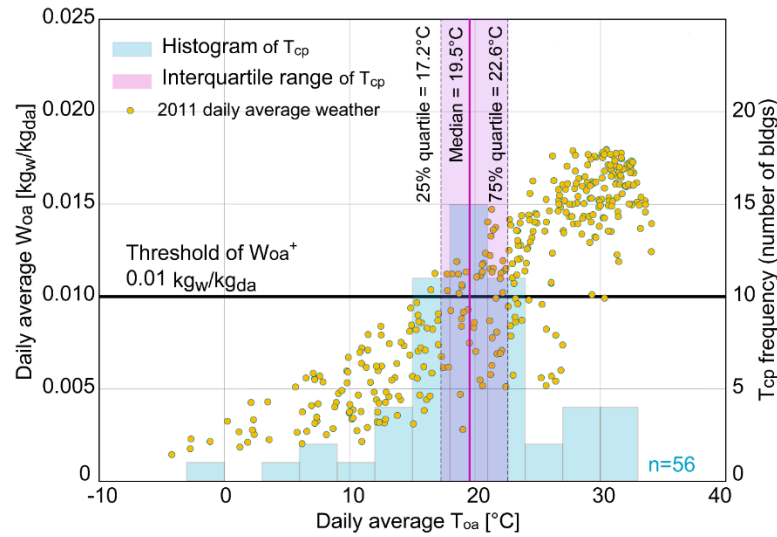


Figure 2-10. The histogram and IQR of T_{cp} in the 4P-CP models for the 56 sample buildings plotted over the weather data: daily average W_{oa} versus daily average T_{oa} .

The IQR of T_{cp} is 17.2–22.6°C (63–73°F). This is higher than the expected cooling coil leaving air temperature T_{cl} for dehumidification around 12.8°C (55°F). The main cause of this discrepancy is the averaging in the daily basis data. During mild

weather with low latent load, dehumidification may occur during only a part of a day. By averaging this transition between dehumidification and non-dehumidification hours for a day, the apparent functional breakpoint of the daily E_{BL} as a function of the daily average T_{oa} becomes higher than T_{cl} under the weather conditions used in this study. The dynamic effects are diminished in daily basis data so that steady-state models can be applied; however, one should acknowledge that the meaning of the T_{cp} estimate from daily basis data is not exactly equal to the T_{cp} in the building energy use model based on the instantaneous energy and mass balances that Katipamula et al. (1994) and Reddy et al. (1995) proposed. In addition to this averaging effect, the reset of the T_{cl} set point during mild weather that does not require the dehumidification and degradation of coil performance may increase the T_{cp} estimates.

The T_{cp} values from the 4P-CP models for some buildings are far from the expected change point, as demonstrated in the histogram in Figure 11. For example, the T_{cp} for building 21 is -0.5°C . This building does not have a clear break point in the typical temperature range for T_{cp} because the effect of W_{oa}^{+} is not significant; however, there are a few points of low E_{BL} data at the lowest temperatures, and the 4P-CP algorithm estimated this as a break point. In this way, the 4P-CP algorithm to find T_{cp} is sensitive to data variations in the tails if the building does not have a significant latent load to create the break point.

2.6 Chapter Summary

This chapter first derived a model structure of steady-state E_{BL} as a linear combination of the T_{oa} , W_{oa}^{+} , E_{sol} , and E_{ele} variables. Based on this structure and a

correlation analysis of the variables, four regression models were proposed, and their applicability was studied using the data from 56 buildings on the Texas A&M University campus. The results demonstrate that the T_{oa} and W_{oa}^+ variables explain the majority of E_{BL} variations and are sufficient for modeling various types of buildings. The mean values of CV-RMSE for the 4P-CP(T), MLR(T,W), MLR(T,W,S), and AR1(T,W) models are 10.4%, 10.0%, 9.7%, and 6.9% respectively. Other factors such as the availability of explanatory variable data, ease of estimation, and ease of parameter interpretation are equally important in practical applications, and these are summarized in Table 2-2.

Table 2-2. Advantages and disadvantages of different models.

Model	4P-CP(T)	MLR(T,W)	MLR(T,W,S)	AR1(T,W)
Required explanatory variables	T_{oa}	T_{oa} and W_{oa}^+	T_{oa} , W_{oa}^+ , and E_{sol}	T_{oa} and W_{oa}^+
Ease of estimation	Easy (simple linear regression with iteration)	Moderate (matrix calculation)	Moderate (matrix calculation)	Complex (matrix calculation and ML estimation)
Ease of parameter interpretation	T_{cp} estimates do not always have physical significance	Easy	Easy	Physical interpretation of autocorrelations is limited.

Even though the estimation of the 4P-CP(T) models requires iteration to find the T_{cp} , it is the simplest method and does not require matrix calculations, so it is widely used for energy use modeling. The residual analysis demonstrates that the 4P-CP(T) models for buildings with a high ventilation rate, such as laboratories and hospitals, have an increasing variance in the high T_{oa} region, and this inhibits its ability to detect metering problems. For these buildings, the MLR models, including the W_{oa}^+ variable, provide superior model fits. The AR1(T,W) models have the smallest residuals among the four models. Despite the good fits, the AR1(T,W) models sometimes capture level changes due to metering problems that should be flagged as outliers. This can be a drawback to the use of this model for data screening.

Overall, the MLR(T,W) model appears to be the most balanced regression model for data screening, considering the availability of the weather variables and the ease of interpretation of the parameter estimates. As the parameters can be directly interpreted in the MLR models, violations of the model assumptions can be detected by a residual analysis. This feature is useful not only to check the validity of the E_{BL} model but also to find the presence of off-assumption building operations such as economizers and temperature set-point resets.

3. ESTIMATION OF BUILDING PARAMETERS USING SIMPLIFIED ENERGY BALANCE MODEL AND METERED WHOLE-BUILDING ENERGY USE²

3.1 Introduction

In this chapter, the E_{BL} models developed in Chapter 2 are used to estimate building parameters. The parameter estimates are compared to actual values to examine the degree of physical significance of the statistically estimated E_{BL} parameters. The use of monthly and daily interval data is evaluated using both synthetic and measured data from three buildings. The methodology follows the studies by Deng (1997) and Reddy et al. (1999), both of which used a variable called the building thermal load Q_B (Reddy et al. 1994). The parameter estimates from the E_{BL} and Q_B variables are compared side by side.

The parameters obtained from data-driven building energy modeling are usually interpreted in terms of the heat-loss coefficient, including ventilation and base temperature, and for dynamic models, thermal capacitance (Sonderegger 1978; Rabl and Rialher 1992; Rabl 1988; Reddy et al. 1999; Sjögren et al. 2009). Hammarsten (1987) and Rabl (1988) noted that there can be considerable uncertainties in the estimated parameters due to non-linearity and non-constancy in actual building systems, measurement errors, and estimation bias from the violation of statistical assumptions.

² Adapted with permission from Masuda, H and Claridge, D. Estimation of building parameters using simplified energy balance model and metered whole building energy use, the 12th International Conference for Enhanced Building Operations. Copyright 2012 the Energy Systems Laboratory, Texas A&M Engineering Experiment Station.

However, several applications of the data-driven method to large numbers of buildings (Sjögren et al. 2007; Liu et al. 2011; Reichmuth and Egnor 2013) demonstrate the usefulness of parameter estimates in building performance benchmarking. Careful interpretation of the parameter estimates may be used for a causal analysis of the changes in the E_{BL} patterns and for a peer-comparison of building load characteristics.

The remainder of this chapter is organized into four sections. Section 3.2 derives mathematical expressions for the key parameters, using regression estimates; Section 3.2 describes the synthetic and actual datasets used in the study; Section 3.4 compares the parameter estimates and the physical values; and Section 0 summarizes the results.

3.2 Formulation of Key Parameters

The key parameters for the E_{BL} and Q_B models are derived based on the assumptions discussed in Section 2.2.1. In addition, the solar gain is assumed to be a linear function of the outside air dry-bulb temperature (Knebel, 1983); and the expression for Q_{sol} is re-written as (Deng 1997):

$$\begin{aligned} Q_{sol} &= a_{sol} + b_{sol} T_{oa} \\ &= a'_{sol} + b_{sol} (T_{oa} - T_{in}) \end{aligned} \quad (3.1)$$

By inserting expressions for Q_{air} , Q_{cond} , and Q_{sol} , given in Equations (2.2), (2.3), and (3.1), into Equation (1.2), the mathematical expressions for the regression parameters in the two-variable E_{BL} model:

$$E_{BL} = \beta_0 + \beta_T T_{oa} + \beta_W W_{oa}^+ + \varepsilon \quad (3.2)$$

can be written as presented in Table 3-1.

The building thermal load Q_B is defined as the cooling energy use minus the heating energy use (Reddy et. al. 1994). Based on the energy balance for the whole building in Equation (1.1), Q_B can be expressed as:

$$\begin{aligned} Q_B &= E_{cool} - E_{heat} \\ &= Q_{air} + Q_{cond} + Q_{sol} + Q_{occ} + Q_{ele}. \end{aligned} \quad (3.3)$$

Deng (1997) and Reddy et al. (1999) introduced two multiplicative correction factors: k_s and k_l . The first factor, k_s , is the ratio of internal sensible loads to measured electricity use for lights and equipment E_{LE} , and the second factor, k_l , is the ratio of the internal latent load to the total internal sensible load, which appears only when latent load exists. If all the internal loads are from the occupants, lights, and equipment, then this relationship is written as:

$$Q_{occ} + Q_{ele} = E_{LE} k_s (1 + k_l X) \quad (3.4)$$

where the indicator variable X is 1 when the latent load exists ($W_{oa} > W_{in}$), and 0 otherwise. Then, the expression of Q_B becomes

$$Q_B = Q_{air} + Q_{cond} + Q_{sol} + E_{LE} k_s (1 + k_l X). \quad (3.5)$$

By inserting expressions for Q_{air} , Q_{cond} , and Q_{sol} from Equations (2.2), (2.3), and (3.1) into Equation (3.5), the MLR for Q_B is expressed as:

$$\begin{aligned} Q_B &= \beta_0 + \beta_{sens} E_{LE} + \beta_{lat} X E_{LE} + \beta_T T_{oa} + \beta_W X (W_{oa} - W_{in}) + \varepsilon \\ &= \beta_0 + \beta_{sens} E_{ele} + \beta_{lat} X E_{LE} + \beta_T T_{oa} + \beta_W W_{oa}^+ + \varepsilon \end{aligned} \quad (3.6)$$

where ε is a random error. The mathematical expressions for the regression parameters can be found in Table 3-1.

Table 3-1. Mathematical expressions for regression model parameters.

Regression parameter	E_{BL}	Q_B
β_0	$(UA_s + m_v c_p + b_{sol})T_{in} - Q_{occ} - a'_{sol}$	$-(UA_s + m_v c_p + b_{sol})T_{in} + a'_{sol}$
β_T	$-(UA_s + m_v c_p + b_{sol})$	$UA_s + m_v c_p + b_{sol}$
β_W	$-m_v h_v$	$m_v h_v$
β_{sens}	Not available	k_s
β_{lat}	Not available	$k_s k_l$

From the regression parameters in Table 3-1, the building parameters and the uncertainties are deduced, as in Table 3-2. The overall heat-loss coefficient estimated from the regression models includes the solar effect. The estimated value, including the solar effect, is designated as U^* , which is defined as $U^* A_s = UA_s + b_{sol}$. The T_{ref} has been defined in Section 2.4.4. This parameter is associated with the T_{in} , and it resembles the balance point temperature (ASHRAE 2009). The physical interpretation of the T_{ref} changes depending on the explanatory variables included in the regression model, which will be discussed later, along with the estimation results.

Table 3-2. Equations to calculate building parameters and the uncertainties from the regression estimates and SEs.

Building parameter	E_{BL}	Q_B	Uncertainty
m_v	$-\hat{\beta}_w / h_v$	$\hat{\beta}_w / h_v$	$\Delta \hat{\beta}_w / h_v$
$U^* A_s$	$-\hat{\beta}_T + \hat{\beta}_w c_p / h_v$	$\hat{\beta}_T - \hat{\beta}_w c_p / h_v$	$\sqrt{(\Delta \hat{\beta}_T)^2 + (\Delta(\hat{\beta}_w \cdot c_p / h_v))^2}$
T_{ref}	$-\hat{\beta}_0 / \hat{\beta}_T$	$-\hat{\beta}_0 / \hat{\beta}_T$	$\left(\frac{\hat{\beta}_0}{\hat{\beta}_T} \right) \sqrt{\left(\frac{\Delta \hat{\beta}_0}{\hat{\beta}_0} \right)^2 + \left(\frac{\Delta \hat{\beta}_T}{\hat{\beta}_T} \right)^2}$
k_s	Not available	$\hat{\beta}_{sens}$	$\Delta \hat{\beta}_{sens}$
k_l	Not available	$\hat{\beta}_{lat} / \hat{\beta}_{sens}$	$\left(\frac{\hat{\beta}_{lat}}{\hat{\beta}_{sens}} \right) \sqrt{\left(\frac{\Delta \hat{\beta}_{lat}}{\hat{\beta}_{lat}} \right)^2 + \left(\frac{\Delta \hat{\beta}_{sens}}{\hat{\beta}_{sens}} \right)^2}$

Note: Δ means an SE, and $\hat{\beta}$ is an estimate of the true parameter value β .

3.3 Data

3.3.1 Synthetic Data

The commercial building reference model (Deru et al. 2011) for EnergyPlus (UIUC and LBNL 2007) simulation software is used to generate synthetic datasets. The existing large office building specifically, which represents construction from 1980 onward in the Houston, TX climate zone, is used. The building has 12 stories above ground, a basement, and a total conditioned area of 46,320 m² (498,588 ft²). Each above-grade floor has six zones: north, east, south, and west perimeters, as well as a core and

plenum. Each floor has a single duct VAV system with reheat terminals, and the building does not use an economizer.

The original input file has three schedule patterns for weekdays (WD), Saturdays (Sat), and Sundays and holidays (Other), as illustrated in Figure 3-1. There are four set points: cooling occupied, heating occupied, cooling unoccupied, and heating unoccupied. During occupied hours, the space temperatures are maintained at the occupied set points. During unoccupied hours, the HVAC systems are disabled until the zone temperatures exceed unoccupied set points. This original input file was used for the As-is case, and another input file was made to generate datasets for two Ideal cases. In the modified input file, the schedules were modified as illustrated in Figure 3-2 so that the parameters can be as constant as possible—ideal for parameter identification. For the Ideal w/ solar case, TMY2 weather data for Houston, TX, was used; meanwhile, for the Ideal w/o solar case, solar insolation values in the TMY2 weather data were set to 0 to remove solar effects from the energy simulation. The combinations of input and weather data files for the three cases can be found in Table 3-3.

Figure 3-3 presents the daily energy uses for electricity (lights, equipment, and fans), cooling, and heating. Figure 3-4 illustrates the E_{BL} and Q_B variables for the plotted As-is case versus the daily average temperature. Similarly, for the Ideal w/o solar case, the daily energy uses are plotted in Figure 3-5, and the E_{BL} and Q_B variables are plotted in Figure 3-6. Note that the signs of the Q_B plots are switched for ease of visual comparison.

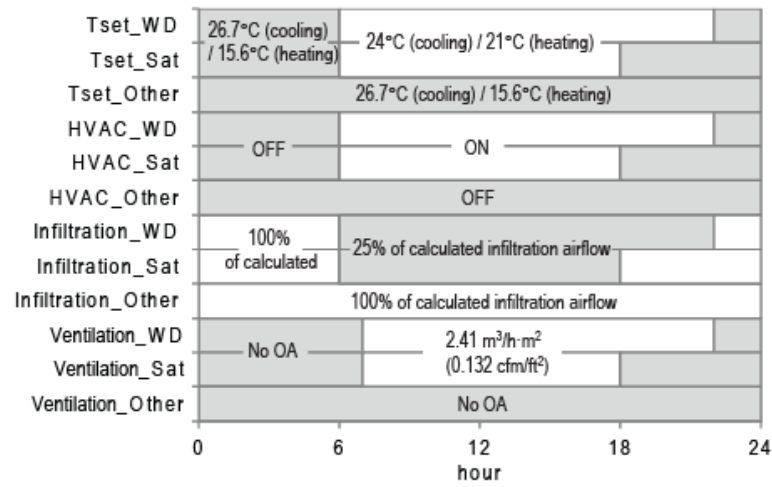


Figure 3-1. System schedules for the As-is case

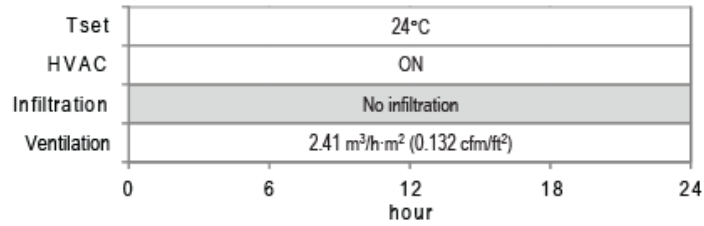


Figure 3-2. System schedules for the Ideal case

Table 3-3. Simulation input conditions for three synthetic datasets

Case designation	System schedules	Weather data
As is	Figure 3-1	TMY2 for Houston, TX
Ideal w/ solar	Figure 3-2	TMY2 for Houston, TX
Ideal w/o solar	Figure 3-2	Modified Houston TMY2 (solar insolation = 0)

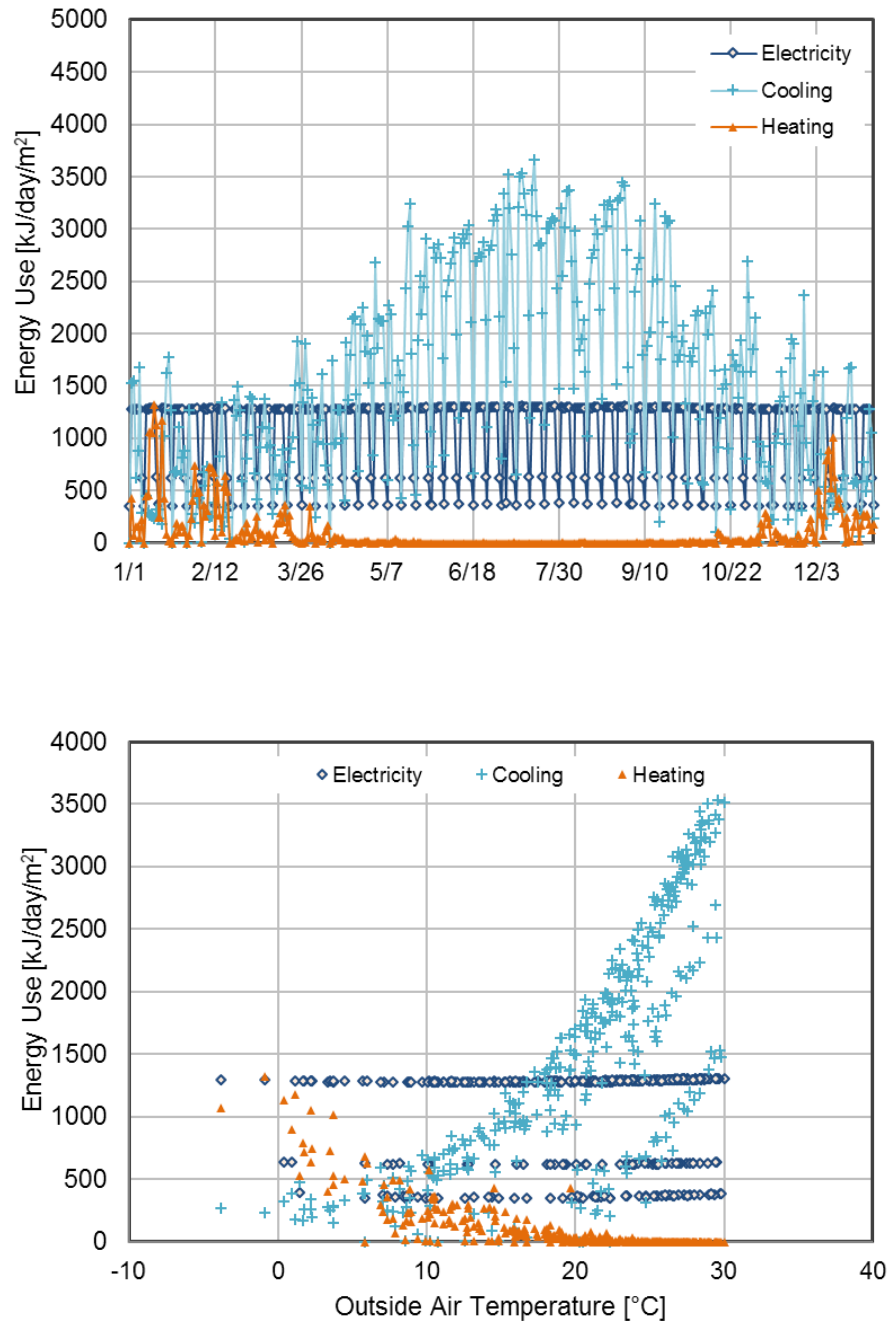


Figure 3-3. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the As-is case. The time series plots are in the top figure and scatter plots versus daily average outside air temperature are in the bottom figure.

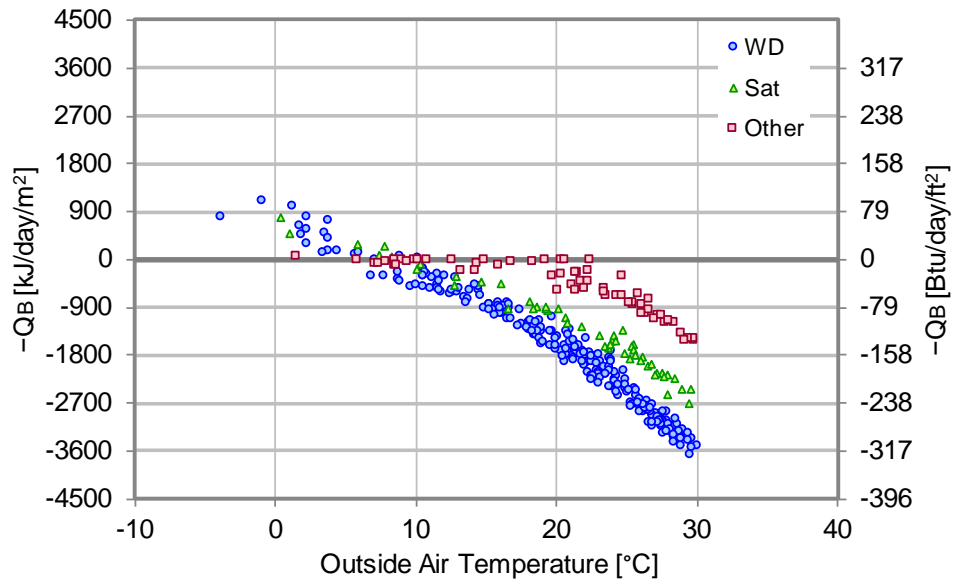
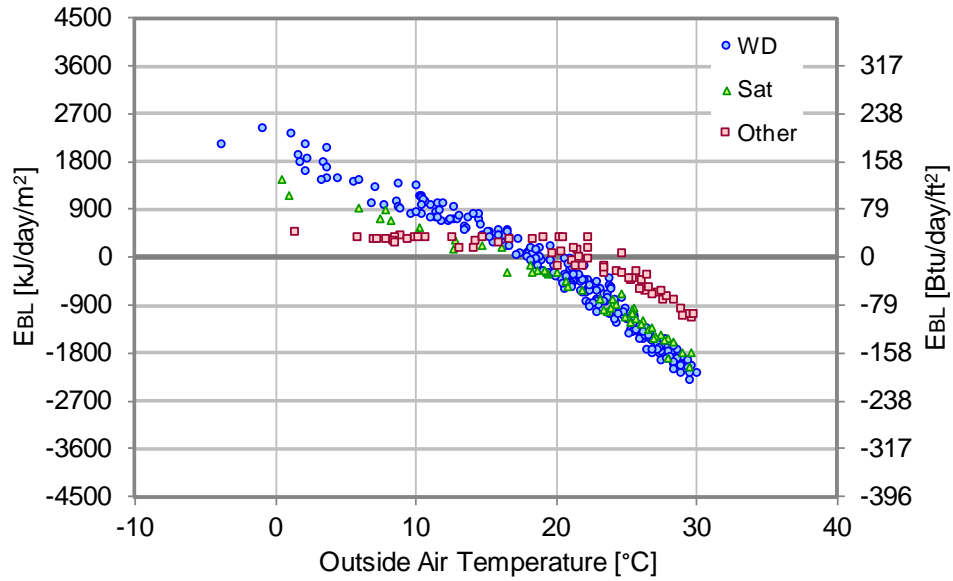


Figure 3-4. E_{BL} and Q_B per unit conditioned floor area for the As-is case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL} .

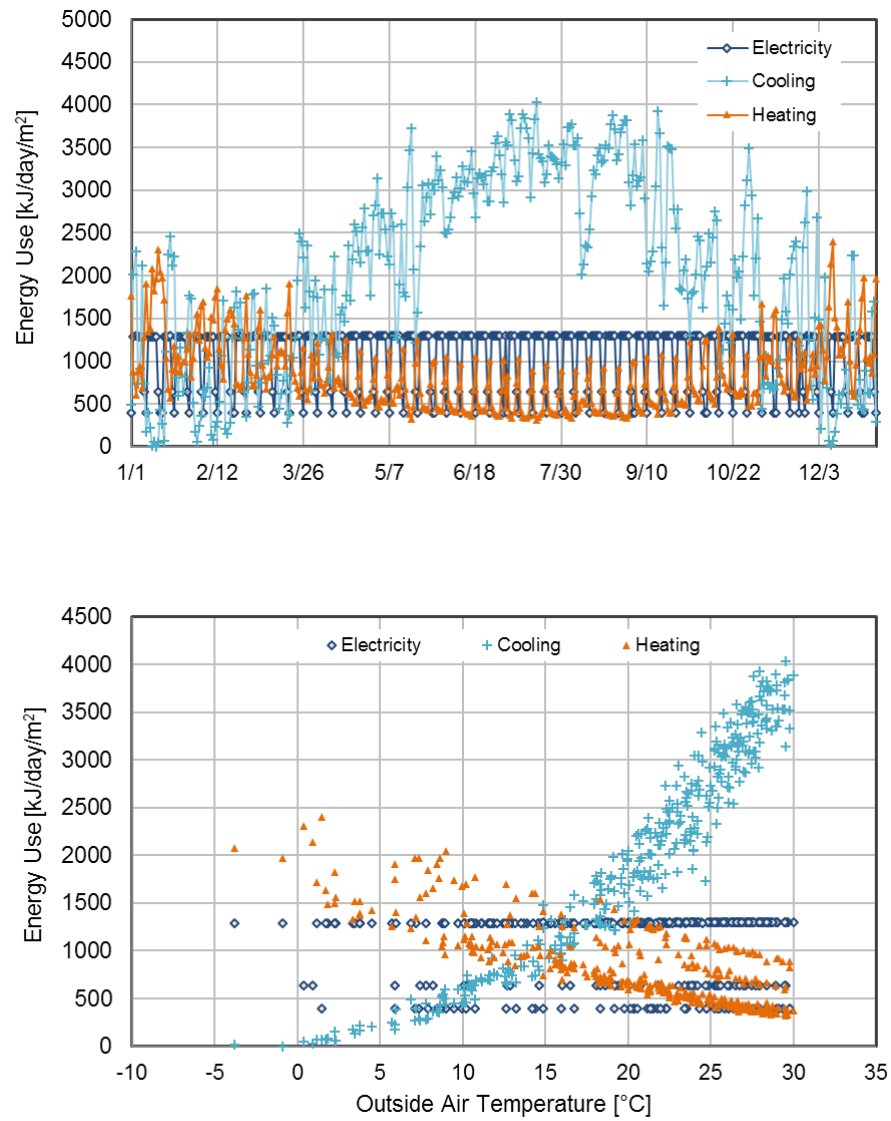


Figure 3-5. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the Ideal w/o solar case. The time series plots are in the top figure and scatter plots versus daily average outside air temperature are in the bottom figure.

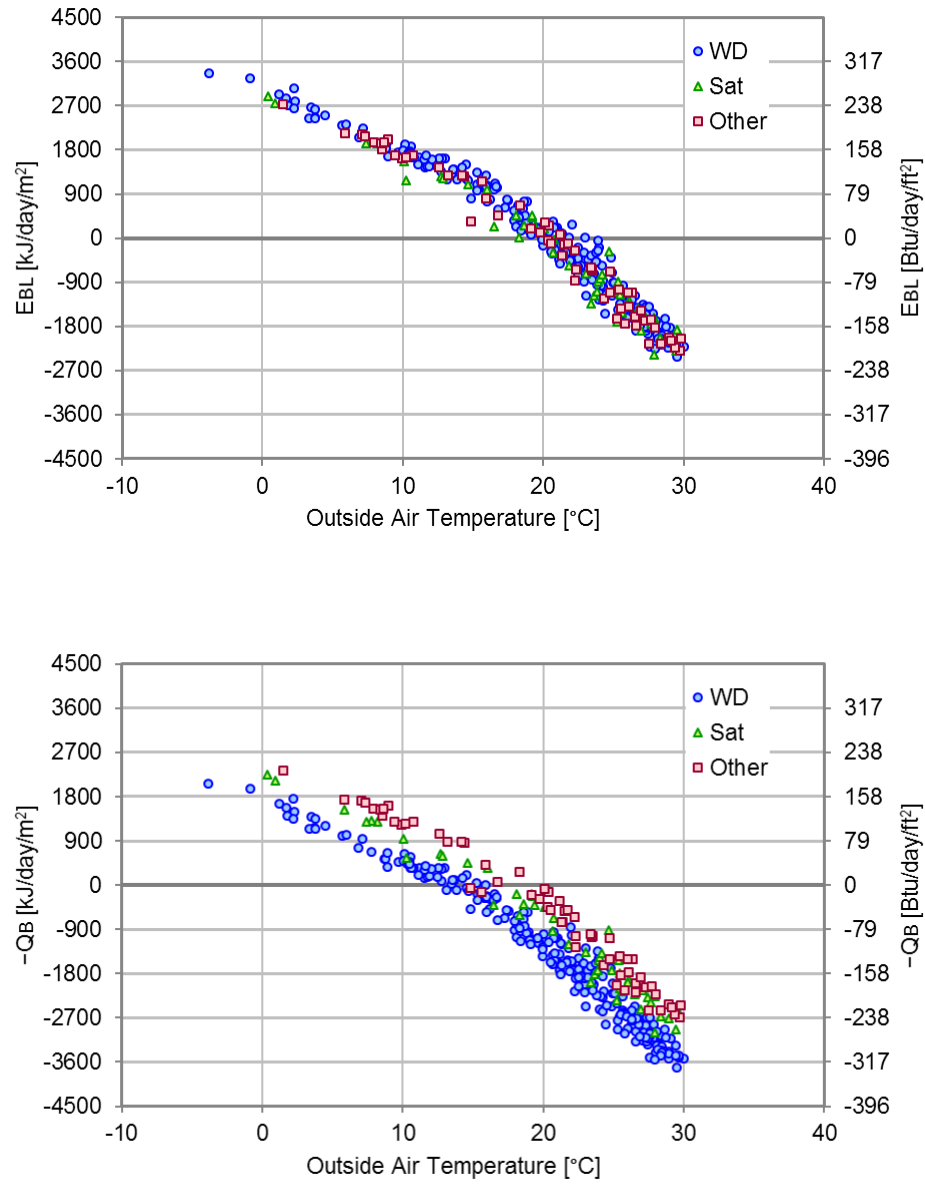


Figure 3-6. E_{BL} and Q_B per unit conditioned floor area in the Ideal w/o solar case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL} .



Figure 3-7. Whole building daily energy uses for electricity, cooling, and heating per unit conditioned floor area for the Ideal w/ solar case. The time series plots are in the top and scatter plots versus daily average outside air temperature is in the bottom figure.

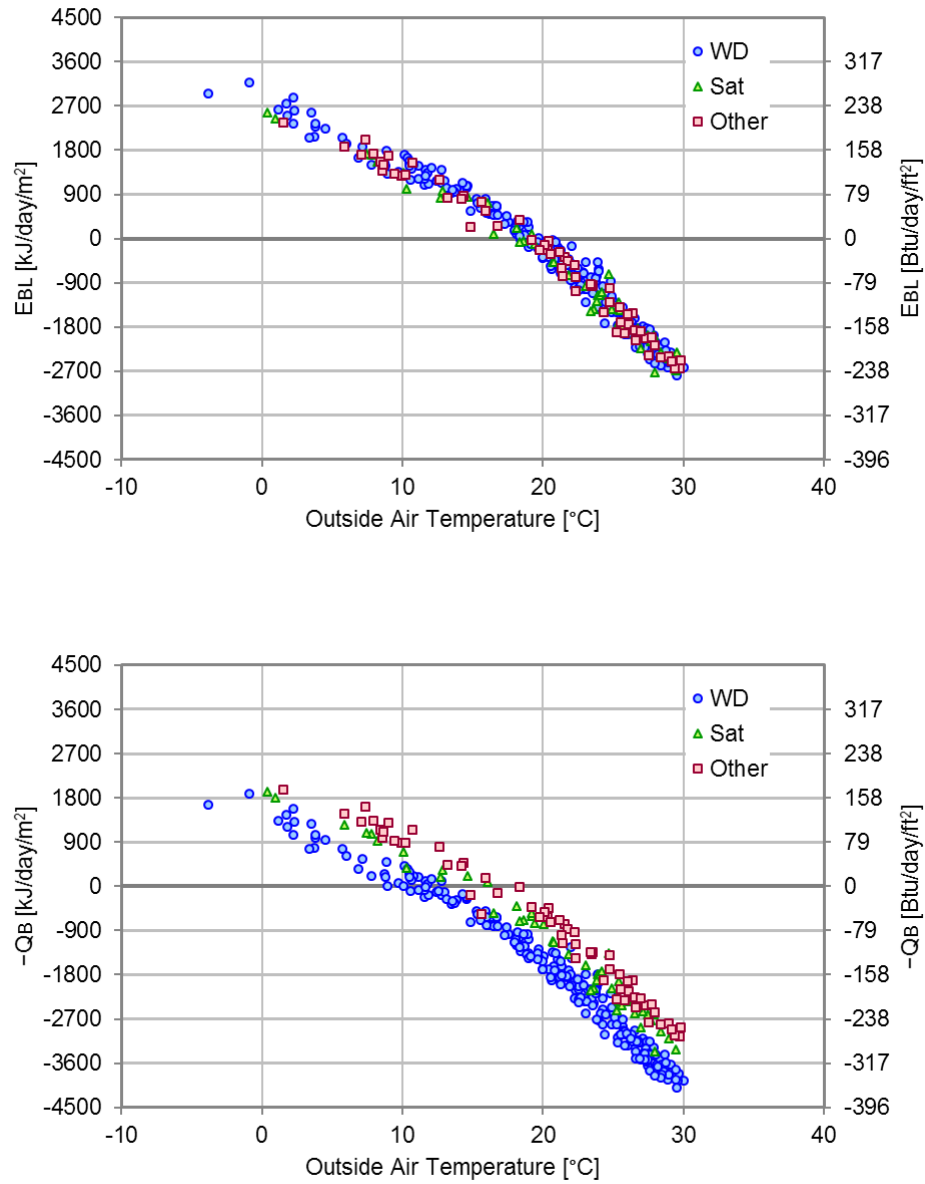


Figure 3-8. E_{BL} and Q_B per unit conditioned floor area in the Ideal w/ solar case plotted versus daily average outside air temperature. The sign of Q_B is flipped for a better comparison with E_{BL} .

3.3.2 Data from Actual Buildings

The whole-building electricity, chilled water, and heating hot water energy use data were collected from the three dormitory buildings that had dedicated outdoor air systems. The HVAC systems were operated continuously in these buildings. The QCLCD for College Station, TX, were obtained from NOAA (2012). The outside airflow rate was measured at the outside air-handling units on April 24, 2012. Furthermore, the measured values of the buildings' total outside airflow rate and the period and the number of days of available E_{BL} and Q_B data are listed in Table 3-4. The daily energy use data are plotted in Figure 3-9, Figure 3-10, and Figure 3-11, and the daily E_{BL} and Q_B data are plotted in Figure 3-12, Figure 3-13, and Figure 3-14, for the Haas, McFadden, and Hobby Hall dormitory buildings respectively.

Table 3-4. Measured outside airflow rates for three dormitory buildings along with the period and the number of days of available E_{BL} and Q_B data.

Bldg. Symbol	Gross floor area	Outside airflow rate measured on April 24, 2012	E_{BL} and Q_B data		
			Available energy use data period	No. of daily data	No. of monthly data
Haas	69,668 ft ² (6,472.4 m ²)	8,779 cfm (14,916 m ³ /h)	July 1, 2011 to June 30, 2012	320	12
McFadden	62,156 ft ² (5,774.5 m ²)	10,025 cfm (17,033 m ³ /h)	September 1, 2011 to June 30, 2012	267	10
Hobby	62,156 ft ² (5,774.5 m ²)	7,750 cfm (13,167 m ³ /h)	July 21, 2011 to June 30, 2012	329	12

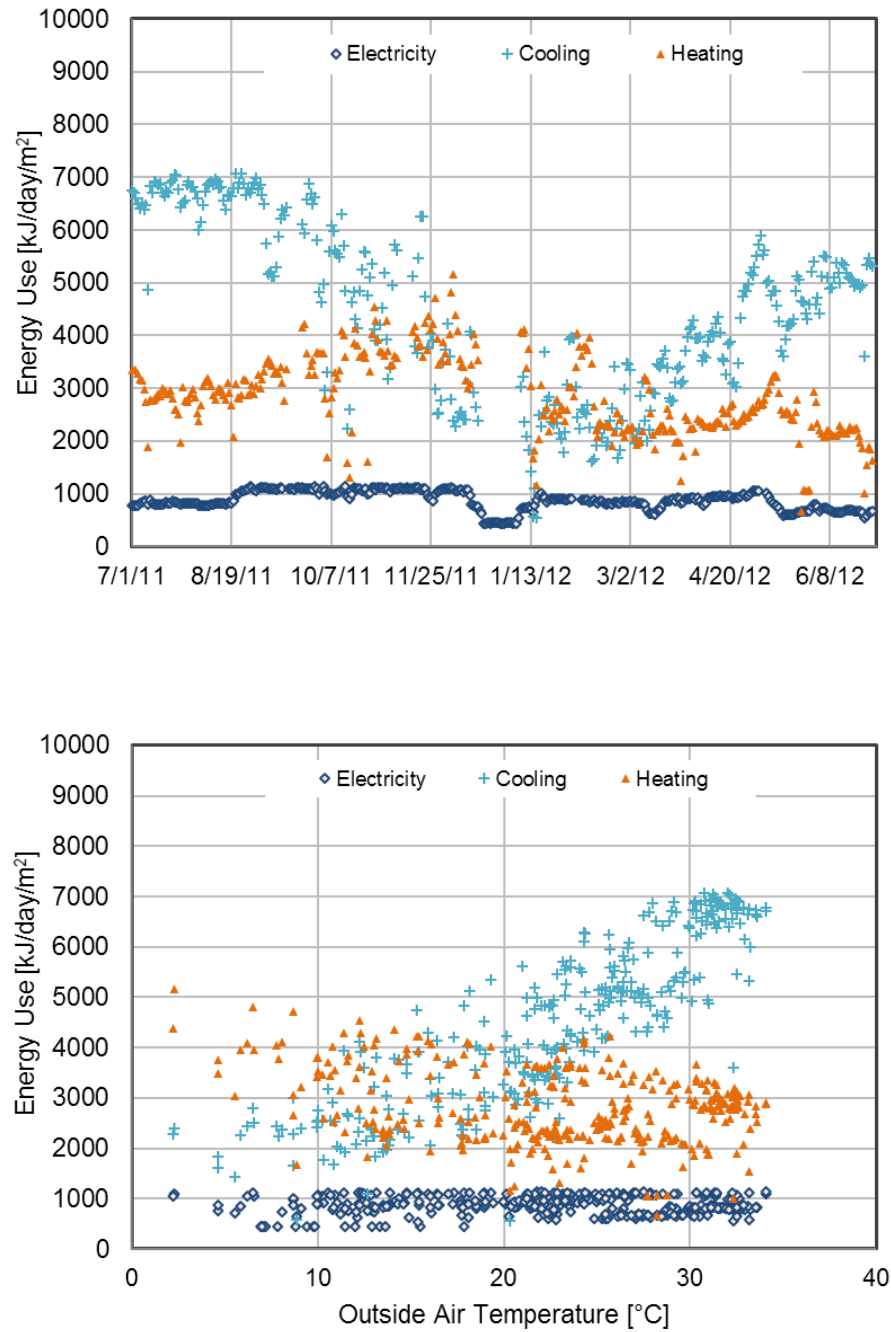


Figure 3-9. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for Haas Hall. Time series plot is on the top and scatter plot versus daily average outside air temperature is on the bottom.

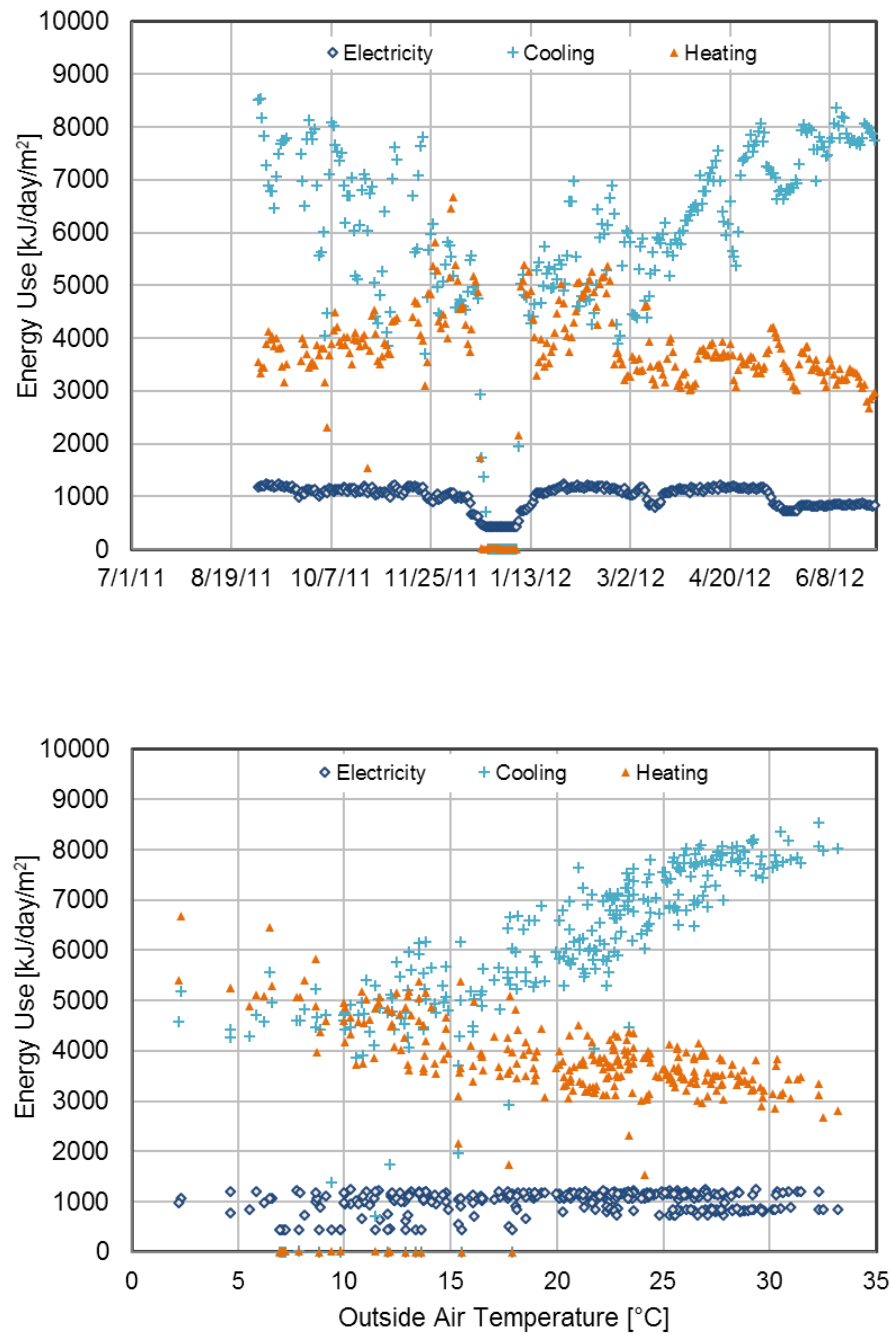


Figure 3-10. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for McFadden Hall. The time series plot is on the top and a scatter plot versus daily average outside air temperature is on the bottom.

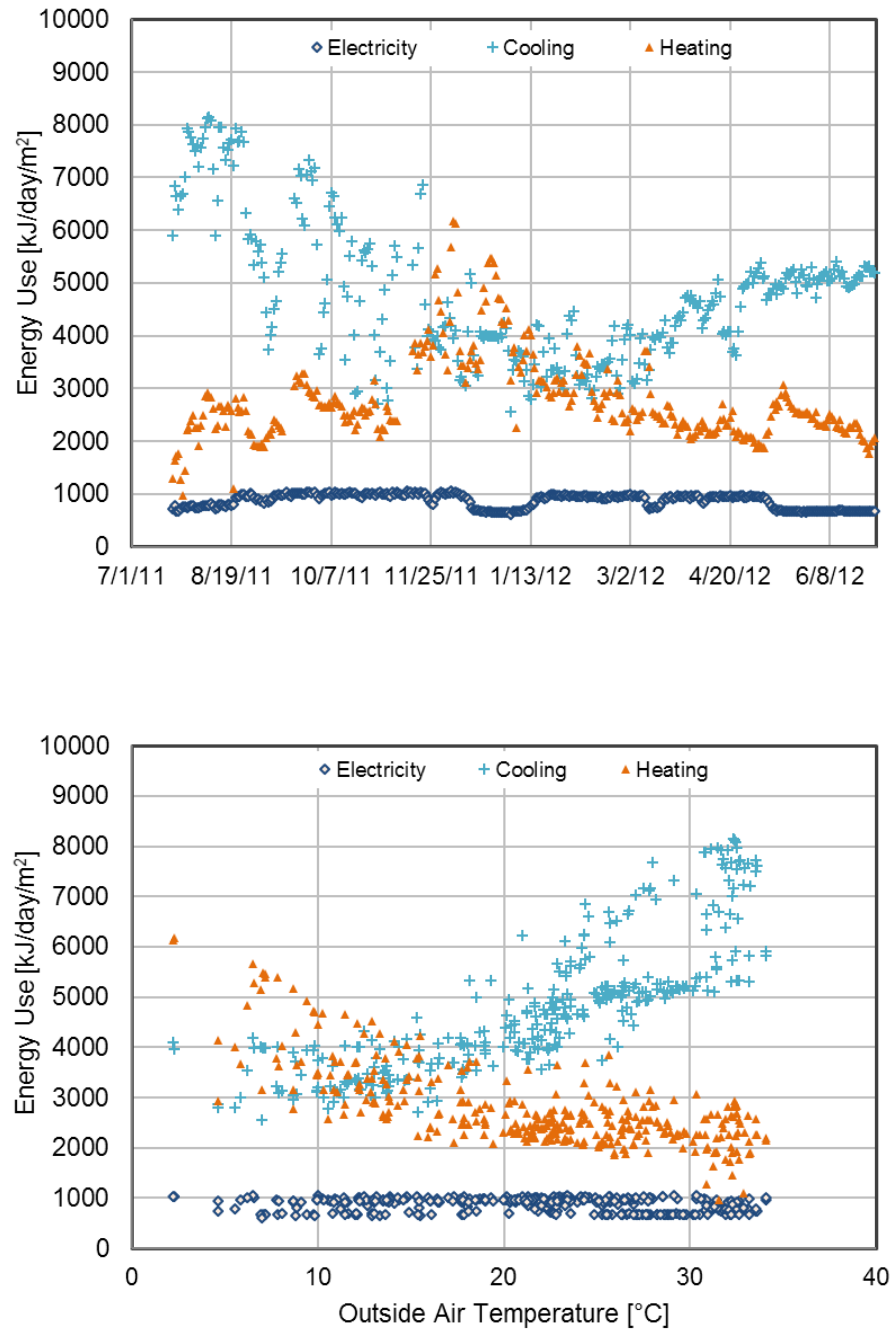


Figure 3-11. Whole building daily energy uses for electricity, cooling, and heating per unit floor area for Hobby Hall. Time series plot is on the top and a scatter plot versus daily average outside air temperature is on the bottom.

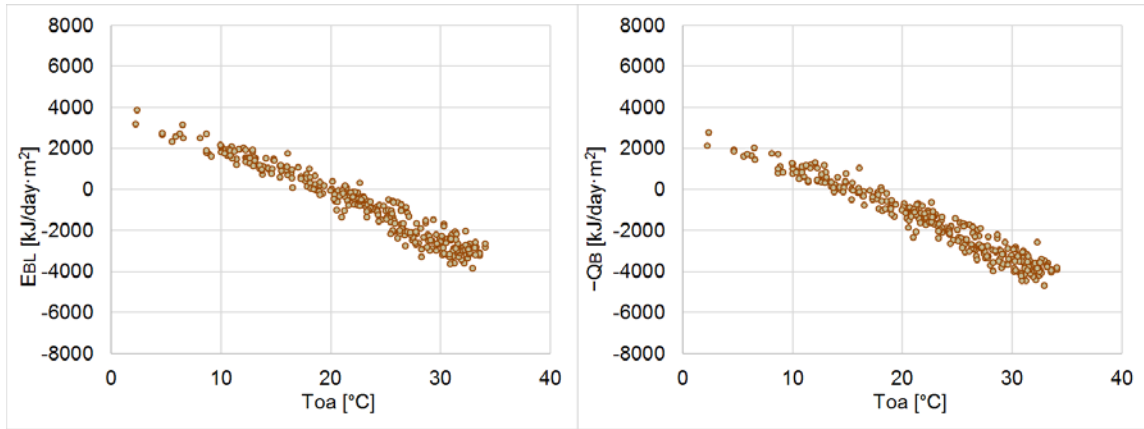


Figure 3-12. E_{BL} and Q_B daily data for Haas Hall during the period July 1, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL} .

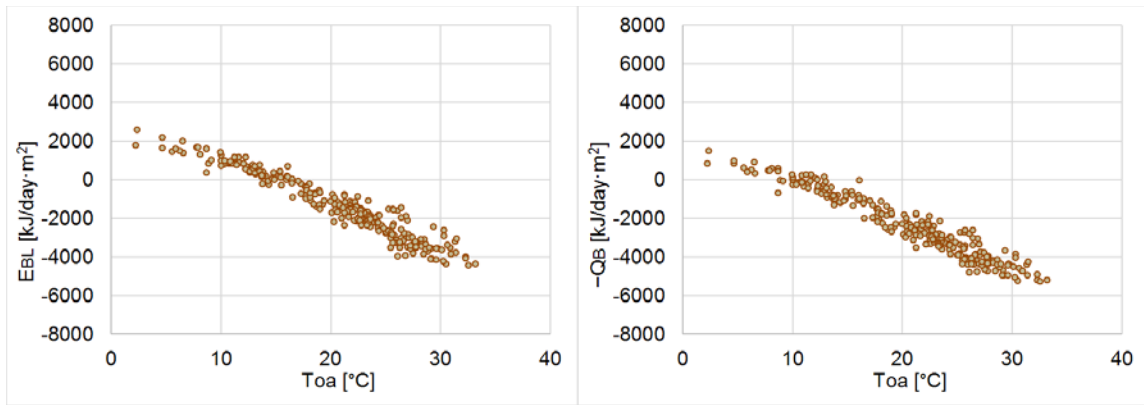


Figure 3-13. E_{BL} and Q_B daily data for McFadden Hall during the period September 1, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL} .

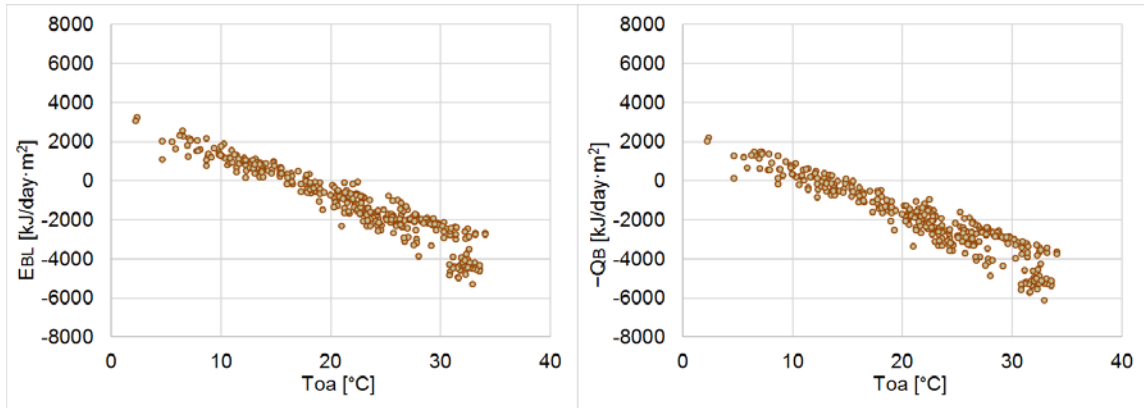


Figure 3-14. The E_{BL} and Q_B daily data for Hobby Hall during the period July 21, 2011–June 30, 2012. The sign of Q_B is flipped for a better comparison with E_{BL} .

3.3.3 Estimation Procedure

The fourteen sets of data listed in Table 3-5 were prepared, and the E_{BL} and Q_B variables were calculated for each set. The data for the As-is case were grouped into three day types—WD, Sat, and Other—when estimated, since the parameters vary between those. The models were estimated with the statistical computing software R (Team 2011). The electricity use variable E_{LE} was not included in the daily interval Q_B models for the As-is cases because the daily electricity use for lights and equipment from the simulation is perfectly constant within the data for each day type, and the parameter estimates become 0. The variable XE_{LE} was removed from the daily Q_B models for the three dormitories and from all the monthly Q_B models because when included, the estimates have reverse signs and/or the estimates are not statistically significant. The explanatory variable terms included in each final model can be found in Table 3-5.

To measure the level of multicollinearity, the Variance Inflation Factors (VIFs) were calculated for each data set. The VIF is defined as (Haan 2002):

$$\text{VIF} = \frac{1}{1 - R_i^2} \quad (3.7)$$

where R_i^2 is the multiple coefficient of determination between the i^{th} explanatory variable and all of the other explanatory variables in the regression equation. The exact value of the VIF at which multicollinearity is declared depends on the individual investigator. Some use a value of 5 and others 10 (Haan 2002).

Table 3-5. Data sets used in the analysis, and the explanatory variable terms included in the regression models. The checked terms are included.

Dataset	Explanatory variable terms included in the regression models					
	E_{BL}		Q_B			
	T_{oa}	W_{oa}^+	T_{oa}	W_{oa}^+	E_{LE}	XE_{LE}
<i>Daily interval</i>						
Ideal w/ solar	✓	✓	✓	✓	✓	✓
Ideal w/o solar	✓	✓	✓	✓	✓	✓
As is (WD)	✓	✓	✓	✓		
As is (Sat)	✓	✓	✓	✓		
As is (Other)	✓	✓	✓	✓		
Haas (Jul–Jun)	✓	✓	✓	✓	✓	
McFadden	✓	✓	✓	✓	✓	
Hobby	✓	✓	✓	✓	✓	
<i>Monthly interval</i>						
Ideal w/ solar	✓	✓	✓	✓	✓	
Ideal w/o solar	✓	✓	✓	✓	✓	
As is	✓	✓	✓	✓	✓	
Haas (Jul–Jun)	✓	✓	✓	✓	✓	
McFadden	✓	✓	✓	✓	✓	
Hobby	✓	✓	✓	✓	✓	

3.4 Results and Discussion

3.4.1 Evaluation Using Synthetic Data

The m_v converted into a volumetric flow rate, the overall heat-loss coefficients U^* , and the absolute value of the temperature slope coefficient $|\beta_T|$, estimated from the daily interval synthetic data, are compared to the assumed true values in Figure 3-15. The signs of the β_T estimates for the E_{BL} and Q_B models are opposite; therefore, the absolute value $|\beta_T|$ is used for comparison. The assumed true values and percent biases for m_v , U^* , and $|\beta_T|$ can be found in Table 3-6, Table 3-7, and Table 3-8 respectively. In these tables, the bias is the difference between the estimate and the assumed true value, and the percent bias means the percentage of the bias relative to the assumed true value. The VIFs of the daily explanatory variables are listed in Table 3-9.

Overall, the E_{BL} and Q_B models have comparable estimates for daily interval synthetic datasets. In the Ideal cases, despite the presence of solar loads, the bias of the m_v estimates using E_{BL} and Q_B models were within 10%. The presence of solar loads increased the percent bias of the $|\beta_T|$ estimates from 4.8% to 13.8% using Q_B models and from 9.0% to 15.8% using E_{BL} models. Similarly, the presence of solar insolation increased the percent bias of the m_v estimates from -2.5% to -9.3% using Q_B models and from 0.03% to -7.9% using E_{BL} models. The U^* estimates involve two regression parameters, namely β_T and β_W , generating larger estimation biases, compared to m_v estimates that involve only one regression parameter: β_W . The bias of the U^* estimates increased from 14.0% to 43.1% using Q_B models and from 20.5% to 45.9% as solar loads were included. It should be noted that the true value for U^* , presented in Figure

3-15 and Table 3-7, is actually for U , which does not include the solar effect, while the U^* estimate includes the solar effect. Therefore, the presented bias is larger than the actual bias.

For the WD and Sat day types in the As-is case, the parameters are estimated fairly well, and they are comparable to the Ideal cases; nevertheless, these simulation models have some exceptions to the model assumptions. The $|\beta_T|$ estimates for the WD and Sat day types are seemingly as good as those for the Ideal cases; however, one should be cautious of this result. The T_{in} decreases with the T_{oa} in the As-is case because of the cooling/heating space temperature set points and system operation schedules, and this can decrease the $|\beta_T|$ estimate. However, the $|\beta_T|$ estimate may already be overestimated due to solar loads. These two factors can balance to result in seemingly good estimates. This type of error can be avoided by using $(T_{oa} - T_{in})$ as a variable instead of using T_{oa} , or by correcting the model using a linear expression for T_{in} as a function of T_{oa} . For the Other day type, meaningful estimates are not available.

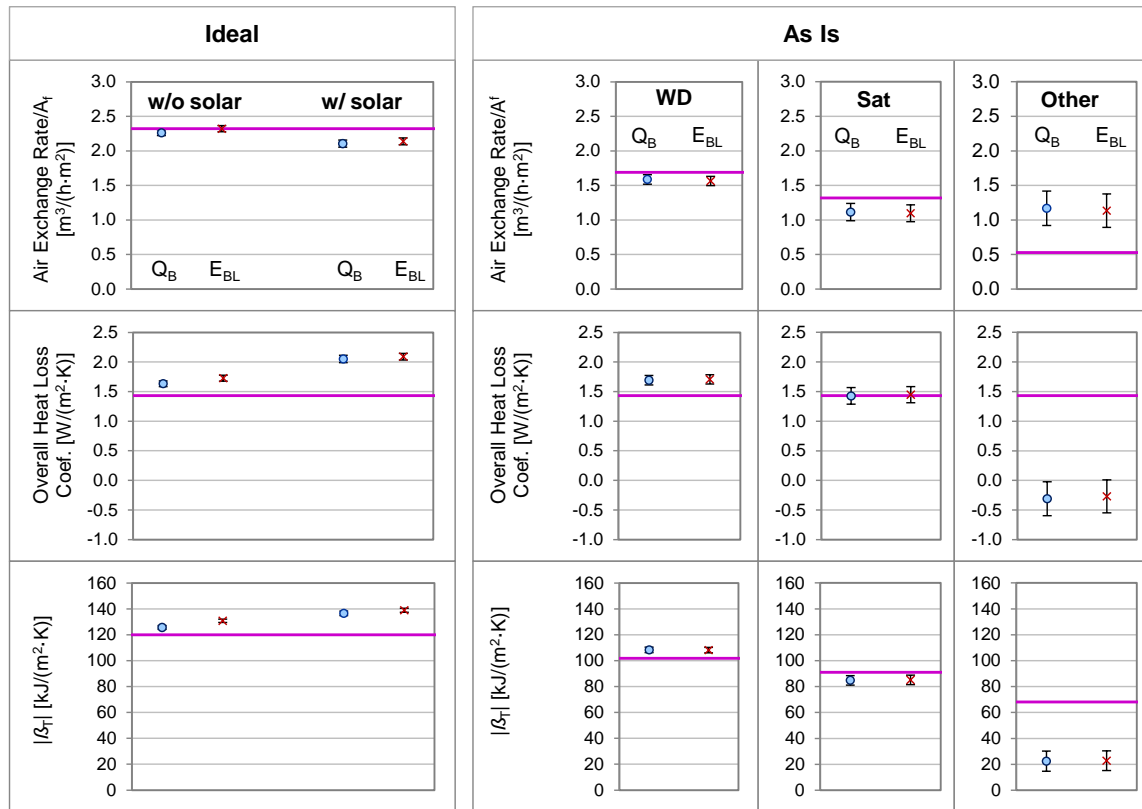


Figure 3-15. Parameter estimates from synthetic daily data for the Ideal and As-is cases. For each of the parameters, the assumed true value is depicted as a solid line, and the parameter estimates using Q_B and E_{BL} are indicated as a circle and a cross respectively, along with the SEs, which are presented as bars.

The parameter estimates using monthly interval synthetic data are presented in Figure 3-16. In the Ideal cases using monthly interval data, the parameter estimates have larger biases, compared to the results from the daily interval data. This may be due to a large collinearity between the explanatory variables in the monthly data, as indicated by the VIFs in Table 3-9. The reason for the good agreement between the estimates and the assumed true values in the As-is case using monthly data is not clear.

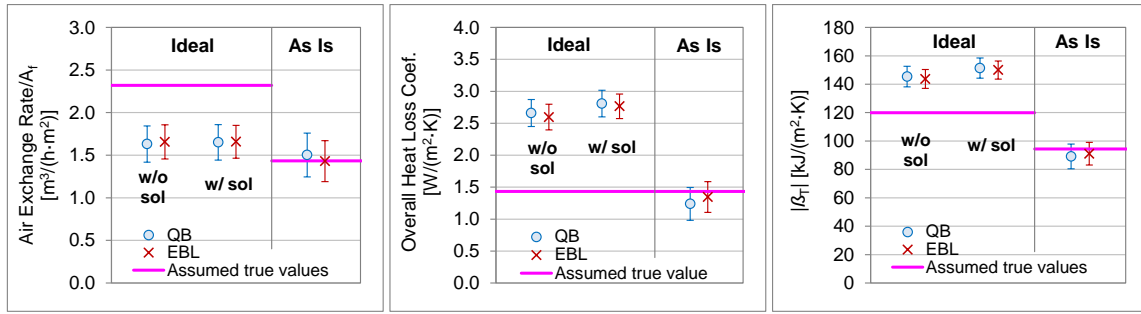


Figure 3-16. Parameter estimates from synthetic monthly data for the Ideal and As-is cases. For each of the parameters, the assumed true value is depicted as a solid line, and the parameter estimates using Q_B and E_{BL} are indicated as a circle and a cross respectively, along with the SEs, which are displayed as bars.

Table 3-6. Assumed true values and percent bias of estimates for m_v .

	True Value [m ³ /(h·m ²)]	Q_B		E_{BL}	
		Daily	Monthly	Daily	Monthly
Ideal w/o solar	2.32	−2.5%	−29.7%	0.03%	−28.6%
Ideal w/ solar	2.32	−9.3%	−28.8%	−7.9%	−28.6%
As-is (WD)	1.69	−6.1%	−11.1%	−7.5%	−15.4%
As-is (Sat)	1.32	−15.5%		−16.8%	
As-is (Other)	0.53	121%		114%	

Table 3-7. Assumed true values and percent bias of estimates for overall heat-loss coefficient U^* .

	True Value [W/(m ² ·K)]	Q_B		E_{BL}	
		Daily	Monthly	Daily	Monthly
Ideal w/o solar	1.43	14.0%	85.7%	20.5%	81.2%
Ideal w/ solar	1.43	43.1%	95.9%	45.9%	93.0%

Table 3-7 Continued.

	True Value [W/(m ² ·K)]	Q_B		E_{BL}	
		Daily	Monthly	Daily	Monthly
As-is (WD)	1.43	18.1%	-13.6%	19.0%	-6.0%
As-is (Sat)	1.43	-0.5%		1.0%	
As-is (Other)	1.43	-122%		-119%	

Table 3-8. Assumed true values and percent bias of estimates for temperature slope $|\beta_T|$.

	True Value [kJ/(m ² ·K)]	Q_B		E_{BL}	
		Daily	Monthly	Daily	Monthly
Ideal w/o solar	119.9	4.8%	21.2%	9.0%	19.8%
Ideal w/ solar	119.9	13.8%	26.2%	15.8%	25.0%
As-is (WD)	101.7	6.5%	-12.4%	6.3%	-10.5%
As-is (Sat)	91.0	-6.8%		-6.4%	
As-is (Other)	68.2	-67.2%		-66.6%	

Table 3-9. The VIFs for explanatory variables in the models for synthetic data. The values for daily (D) and monthly (M) data are compared.

Explanatory variable	D	M	D	M	D	M	D WD	D Sat	D Other
T_{oa}	4.02	15.63	3.07	8.34	3.07	7.59	2.96	2.95	3.77
W_{oa}^+	3.15	8.12	3.08	7.94	3.07	7.59	2.96	2.95	3.77
E_{LE}	1.55	1.28	1.00	1.14					
XE_{LE}	2.97	5.97							

The T_{ref} estimates for the synthetic data sets are listed in Figure 3-17 along with the distribution of the daily average T_{in} in the building. Furthermore, T_{in} is plotted with T_{ref} because T_{in} has a large influence on T_{ref} . The physical significance of T_{ref} varies with the structure of the regression models—the different mathematical expressions and approximate values can be found in Table 3-10. Both the E_{BL} and Q_B models have good estimates for T_{ref} in the Ideal cases. In the As-is case, the bias and estimation errors increase as the unconditioned hours increase. The T_{ref} estimates from the E_{BL} models appear to be more stable across the different data sets if the HVAC systems run for at least 16 hours per day.

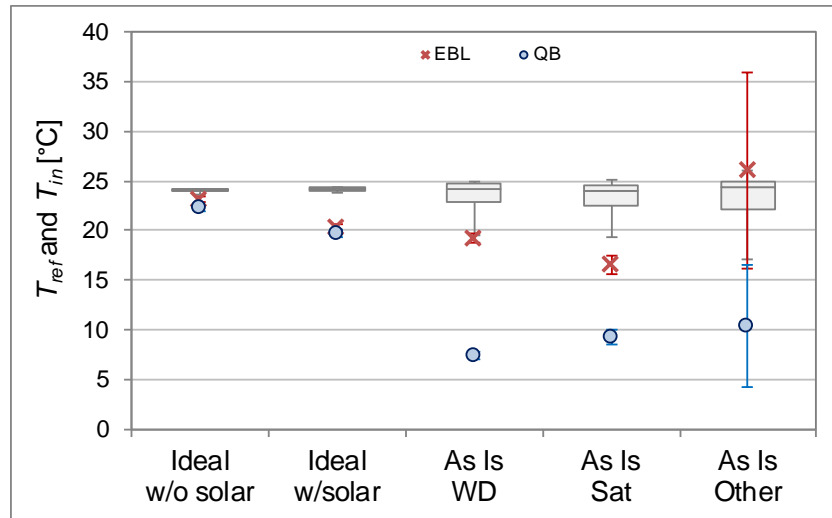


Figure 3-17. The T_{ref} estimates and the distributions of T_{in} . For each case, estimates using E_{BL} and Q_B are depicted with the SEs, and the annual distribution of the daily average T_{in} is presented by box-and-whisker plots.

Table 3-10. Physical meaning of the reference parameter T_{ref} for different models used for synthetic data. Expected values are given as well.

Case	E_{BL}	Q_B
Ideal w/o sol	$T_{in} - \frac{Q_{occ}}{UA_s + m_v c_p}$ <p>~23.2°C (73.8°F)</p>	T_{in} <p>~24.0°C (75.2°F)</p>
Ideal w/ sol	$T_{in} - \frac{Q_{occ} + a'_{sol}}{UA_s + m_v c_p + b_{sol}}$ <p>~21.8°C (71.2°F)</p>	$T_{in} - \frac{a'_{sol}}{UA_s + m_v c_p + b_{sol}}$ <p>~22.6°C (72.7°F)</p>
As is	$T_{in} - \frac{Q_{occ} + a'_{sol}}{UA_s + m_v c_p + b_{sol}}$ <p>~21.8°C (71.2°F)</p>	$T_{in} - \frac{Q_{occ} + a'_{sol} + E_{ele}}{UA_s + m_v c_p + b_{sol}}$ <p>~13.8°C (56.8°F)</p>

3.4.2 Application to the Data from Actual Buildings

The outside airflow estimates from the daily and monthly interval data for three dormitory buildings are compared to the measured values in Figure 3-18, Figure 3-19, and Figure 3-20. The T_{ref} values for the corresponding data are presented in these figures as well. The measured values and percent biases can be found in Table 3-11, and the VIFs of the explanatory variables can be found in Table 3-12.

Overall, the estimates using daily data had similar values for the Q_B and E_{BL} models. However, with monthly data, the estimates using E_{BL} models gave better results. The estimates from the monthly Q_B models have larger SEs, compared to monthly E_{BL}

models, indicating an unstable estimation performance of Q_B models using monthly data. The Q_B model has an extra variable, E_{LE} , which increases the collinearity in the model, especially for the monthly data. The VIFs of the explanatory variables included in the monthly Q_B models for Haas and Hobby Halls are alarmingly high (> 10). As discussed, the presence of a severe collinearity makes the model unreliable. The estimate for Hobby Hall using the monthly Q_B model has a 140% bias with a large SE. In fact, the effect of the E_{EL} variable is overestimated at approximately five times as large as the values for the other buildings. Another problem with Hobby Hall is that the consumption pattern changed during the data period. This can be clearly seen in the scatter plots in Figure 3-11 and Figure 3-14, and it is highly possible that the m_v for this building changed during the data period because of the visible change in the T_{oa} slope. This does not satisfy one of the assumptions of the regression models, namely that parameters are constant during the data period, and it causes larger bias in the estimates for Hobby Hall. The abnormal estimates are indicated by the high value of T_{ref} , near 50°C , for the Hobby Hall monthly Q_B model; this is not a realistic value, based on the physical significance of the parameter.

McFadden Hall has comparable results between daily and monthly data, unlike the other buildings. The VIFs of the monthly data for McFadden Hall are small (< 6), compared to the other buildings, due to a lack of data for the hot and humid months of July and August. This lesser collinearity in the monthly data might be the reason for the similar results between the daily and monthly data.

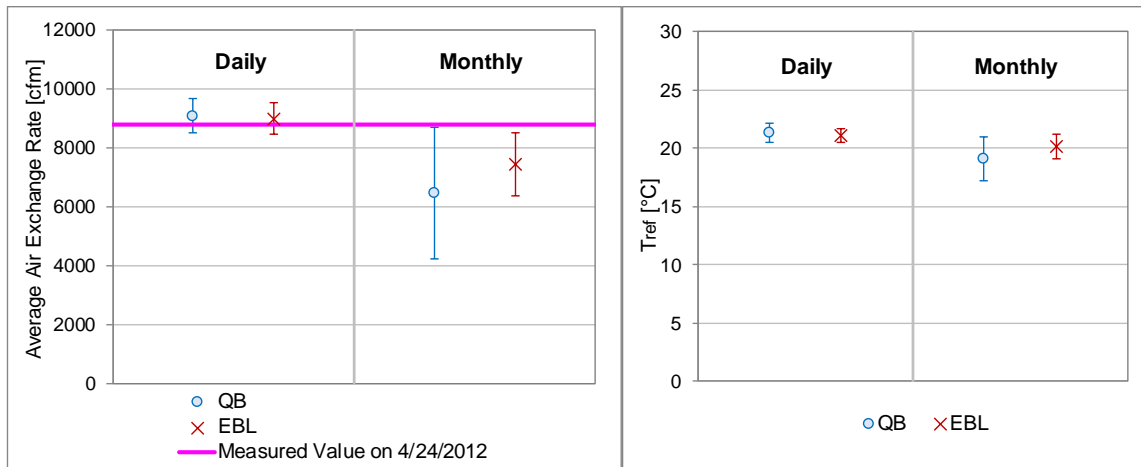


Figure 3-18. Daily average outside airflow rate (left) and T_{ref} (right) estimated for Haas Hall, comparing the estimates using daily and monthly interval data. Two different data periods are used. The SE is presented with bars for each estimate. 1 cfm = 1.699 m³/h.

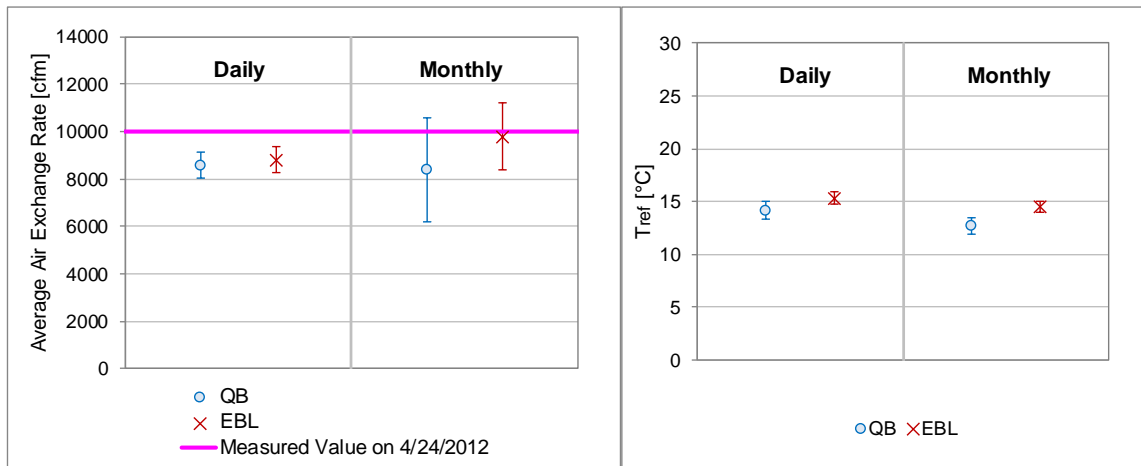


Figure 3-19. Daily average outside airflow rate (left) and T_{ref} (right) estimated for McFadden Hall, comparing the estimates using daily and monthly interval data. The SE is depicted with bars for each estimate. 1 cfm = 1.699 m³/h.

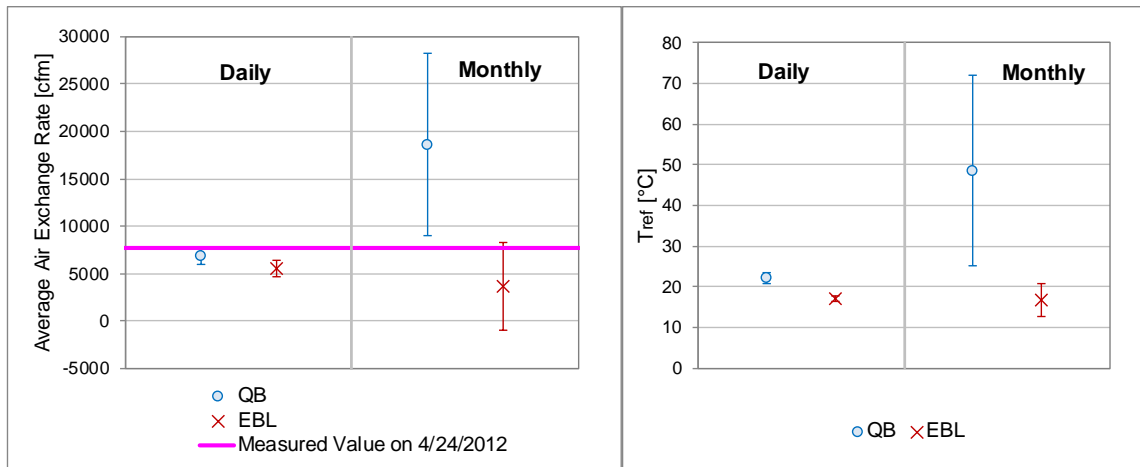


Figure 3-20. Daily average outside air flow rate (left) and T_{ref} (right) estimated for Hobby Hall, comparing the estimates using daily and monthly interval data. The SE is displayed with bars for each estimate. 1 cfm = 1.699 m³/h.

Table 3-11. True values and bias of the m_v estimates for three dormitory buildings. The bias is expressed as a percentage of the measured value.

Building Name	Measured Value (cfm)	Q_B		E_{BL}	
		Daily	Monthly	Daily	Monthly
Haas	8,779	3.4%	-26.3%	2.5%	-15.4%
McFadden	10,025	-14.6%	-16.5%	-12.0%	-2.3%
Hobby	7,750	-12.2%	140%	-27.8%	-52.6%

Table 3-12. The VIFs for explanatory variables in the models for three dormitory buildings. The values for daily (D) and monthly (M) data are compared.

Explanatory variable	Haas				McFadden				Hobby			
	D	M	D	M	D	M	D	M	D	M	D	M
T_{oa}	2.95	14.5	2.86	4.25	2.17	5.23	2.13	2.48	2.72	11.4	2.66	4.01
W_{oa}^+	3.16	15.3	2.86	4.25	2.33	5.52	2.13	2.48	2.96	21.5	2.66	4.01
E_{LE}	1.13	3.69			1.11	2.33			1.15	6.19		

3.5 Chapter Summary

This chapter examined the degree of physical significance, specifically about the m_v , and overall heat exchange coefficient U^* , of statistically estimated building parameters using E_{BL} and Q_B models. The T_{ref} proposed in 2 was formulated and estimated as well. Estimates using daily data generally demonstrate better accuracy, compared to estimates using monthly data. In synthetic data based on the commercial reference building model, the biases of the m_v estimates were within 10%, and the biases of the temperature slope $|\beta_T|$ estimates were within 16%, using daily Q_B and E_{BL} models, if the HVAC systems operate for 12 hours a day or longer (WD and Sat schedules). The estimation of U^* involves two regression parameters, namely β_T and β_W , and the estimation bias is larger, compared to the m_v estimate that involves one regression parameter: β_W . The biases of the U^* estimates using daily data were within 46%. As expected from the previous studies, the presence of solar loads increases the bias of parameter estimates. In synthetic data, the bias of the m_v estimates increased from 0% to -8%, and the bias of the overall heat-loss coefficient increased from 21% to 46%, using daily E_{BL} models, when the solar insolation is included in the weather data for simulation. The biases of the m_v estimates for three actual buildings using daily E_{BL} and Q_B models were within 28%.

The definition of the Q_B variable resembles the one for E_{BL} , and the E_{BL} and Q_B models generally have a similar degree of accuracy in the parameter estimation. However, the effects of the electricity variables E_{EL} and XE_{EL} in the Q_B models were not estimated properly using monthly data because additional explanatory variables— E_{EL}

and XE_{EL} —increased the collinearity. The estimate for Hobby Hall using the monthly Q_B model had a 140% bias with large uncertainty due to high collinearity. The sample size in monthly data is generally small, and the inclusion of multiple variables should be avoided to decrease the chance of overfitting that can cause misleading estimates.

The T_{ref} has physical significance in temperature that can be easily interpreted. Therefore, it can be utilized as a means to detect physically impossible estimates due to model misspecification or metering problems. The erratic m_v for Hobby Hall using the monthly Q_B model was indicated by a T_{ref} of 50°C, where T_{ref} is supposed to be lower than the space temperature T_{in} .

In the data used in this chapter, m_v had consistently reasonable estimates. The estimation of m_v depends on the outdoor air humidity ratio variable, and if the data lacks hot and humid ambient conditions, then the estimates may not be reliable. Not only can this be caused by missing data, but it can also result from the dry climate where the building stands. The applicability of the method to the different climate zones should be scrutinized in future studies.

4. THE VARIATION OF THE INDOOR REFERENCE VARIABLE

4.1 Introduction

In Chapter 2, a variable called the indoor reference temperature T_{ref} was introduced. This variable is an estimate of the temperature where $E_{BL}=0$ when the E_{BL} versus T_{oa} is plotted. This temperature has often been used as an index to find some types of metering problems such as a calculation error in the thermal energy meter. For instance, in the data presented in Figure 4-1, there is a sudden change in the chilled water and heating hot water energy consumption, resulting in two distinct E_{BL} patterns. In terms of the group of data—before or after the change—that is more likely to be correct, the temperature at $E_{BL}=0$ indicates that the data after the change may be more realistic because it is near the space temperature.

The mathematical expression of T_{ref} in Equation (2.11) demonstrates that the T_{ref} value is T_{in} minus a term that is a function of Q_{occ} , Q_{sol} , E_{ele} , m_v , and UA_s . Although a range, such as 60–80°F, has been used as a rule of thumb for T_{ref} , there have not been any studies on the possible range for T_{ref} in different types of buildings. In this chapter, a possible range for T_{ref} in different types of buildings on the Texas A&M University campus will be estimated. A statistical factor analysis is conducted using simulated data.

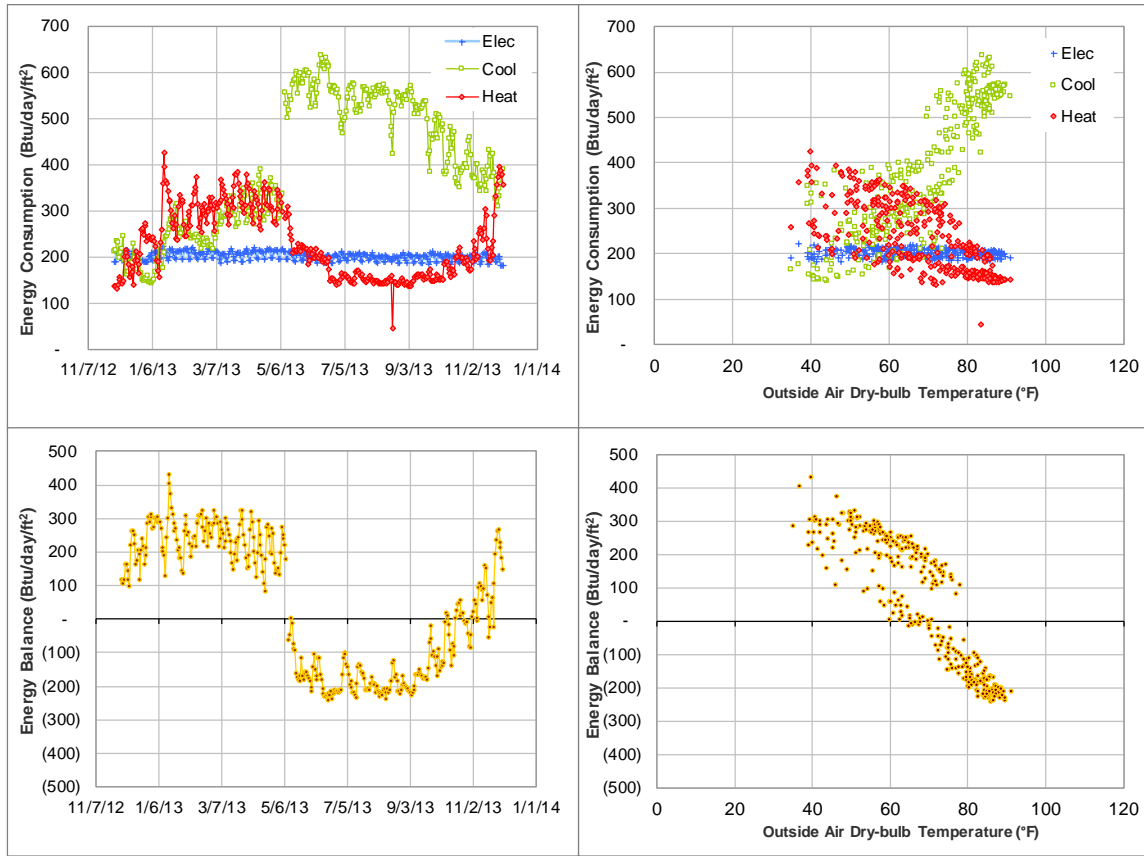


Figure 4-1. Daily energy use and E_{BL} for Heep Laboratory between December 1, 2012 and December 31, 2013.

4.2 Methodology

4.2.1 Sensitivity and Uncertainty Analysis

Saltelli et al. (2008) defines a *sensitivity analysis* as the study of how uncertainty in the output of a model (numerical or otherwise) can be apportioned to different sources of uncertainty in the model input. An *uncertainty analysis* is distinguished from a *sensitivity analysis* when the focus is on quantifying uncertainty in model output (Saltelli et al. 2008). A considerable number of sensitivity analysis methods have been proposed in the literature, and a broad review of sensitivity analysis methods is presented in

Hamby (1994). In this study, regression analysis using a stepwise regression procedure, which is a highly comprehensive technique that is easy to perform with commercially available software (Hamby 1994), was used. First, the factorial design was used to run multiple simulations using eQuest® (LBNL 2011) to generate multiple sets of inputs and outputs for the estimation of a simplified response surface for T_{ref} . The response surface is a regression equation that approximates model output using only the most sensitive model input parameters. Then, a random sampling of the input parameters was conducted to assess the combined variability in the response resulting from considering input parameters simultaneously.

The multiple regression model for E_{BL} with the T_{oa} and W_{oa}^+ variables, as in Equation (2.6), designated as MLR(T,W) in Chapter 2, was used for the estimation of T_{ref} , which was estimated from the regression parameter estimates using Equation (2.11).

4.2.2 Factorial Design

Two earlier studies have performed sensitivity and uncertainty analyses on E_{BL} . Shao (2006) estimates the uncertainty in simplified spreadsheet simulations due to the variations of input parameters using a factor analysis. In this factor analysis, annual RMSE comparing the simulated data with the measured data, is analyzed. Seven factors—the solar heat gain coefficient F , the U-value of windows U_{window} , the U-value of walls U_{wall} , the indoor air temperature T_R , the outside air intake rate V_{OA} , the cooling coil leaving air temperature T_{CL} , and the heat load due to occupants Q_{occ} —are included in the 2^{7-3} fractional factorial design. The factor levels used by Shao can be found in Table 4-1. The most influential factors were determined to be V_{OA} , T_{CL} , and T_R .

Table 4-1. Factors and levels used in Shao (2006).

	Parameter	-1	+1	Unit
Factor 1	F	0.25	0.87	
Factor 2	U_{window}	0.1	1.04	$Btu / hr \cdot ft^2 \cdot ^\circ F$
Factor 3	U_{wall}	0.1	0.2	$Btu / hr \cdot ft^2 \cdot ^\circ F$
Factor 4	T_R	65	80	$^\circ F$
Factor 5	V_{OA}	0.05	0.8	cfm / ft^2
Factor 6	T_{CL}	50	70	$^\circ F$
Factor 7	Q_{occ}	3	8	$Btu / ft^2 \cdot day$

Ji et al. (2008) conducted factor analyses to investigate the effects of factors on the slopes and intercepts for the E_{BL} vs. T_{oa} and E_{BL} vs. h_{oa} simple linear regression models, using simplified spreadsheet simulations. Five factors, namely the OA ratio, the T_{cl} , zone temperature T_z , UA value, and occupant density are included in the 2^5 full factorial design. The factor levels used in Ji et al. can be found in Table 4-2.

Table 4-2. Factors and levels used in Ji et al. (2008).

	-1	+1
OA Ratio	5%	15%
T_{CL}	45°F	65°F
T_z	68°F	82°F
UA	5000 Btu/(hr*F)	15000 Btu/(hr*F)
Occupant density	300 ft ² /person	500 ft ² /person

The result of Ji et al. demonstrated that the OA ratio and UA value have by far the strongest positive effects on the intercept, followed by cold deck temperature in both the E_{BL} vs. T_{oa} and the E_{BL} vs. h_{oa} models. Similarly, the OA ratio and UA value have strong negative effects on the slope, followed by the zone temperature in both the E_{BL} vs. T_{oa}

and the E_{BL} vs. h_{oa} models, whereas the cold deck temperature has a strong positive effect on the slope.

These studies demonstrated a high significance of factors such as the outside air flow, the overall heat transfer coefficient of building envelope, the indoor temperature, and the cooling coil leaving air temperature.

Based on the results from previous studies and the mathematical expression of T_{ref} in Equation (2.11), 11 factors were included in the factor analyses in this study. They are the aspect ratio of building floors, the window glass types, the cold deck temperature, the building direction, the number of floors, the opaque surface absorptance, the indoor temperature, the occupancy density, the outside air flow, the opaque surface U-values, and the window ratio. While Q_{sol} and UA_s depend on the building surface area A_s , Q_{occ} and m_v usually increase with the building volume or floor area A_f . Therefore, T_{ref} and β_T may be affected by the surface area to floor area ratio A_s/A_f . The aspect ratio of floors and the number of floors are included to observe the effect of A_s/A_f on each parameter.

To assign factor levels effectively with a small number of simulation runs, definitive screening design (DSD), which is a new class of design proposed by Jones and Nachtsheim (2011), was used. For m factors with three levels, DSDs require only $2m+1$ runs. Since there are three-level designs, curvature in the relationship between any factor and the response of interest can be analyzed. The main effects are completely independent of two-factor interactions, and none of the two-factor interactions are confounded with any other two-factor interactions. Therefore, there is no need to

conduct additional analyses to identify the effect from confounding factors in the first analysis.

The levels for the four categorical factors and seven continuous factors are defined in Table 4-3 and Table 4-4 respectively. To easily change the eQuest simulation inputs, the aspect ratio, glass types, and cold deck temperature are treated as categorical factors. The minimum and maximum values of the factor levels are determined so they roughly represent the ranges of values for buildings on the Texas A&M University campus. These factor levels are assigned to 26 simulation runs using DSD, as indicated in Table 4-5.

Table 4-3. Factors and levels (categorical factors).



Categorical Factors	Level 1	Level 2
Aspect ratio of floors	1:1 (square) 	1:4 
Glass types	Single Tint Grey	Double Low-E (e2 = .1) Clear
Cold deck temperature reset	55°F const.	55°F–65°F reset

Table 4-4. Factors and levels (continuous factors).




Continuous Factors	-1	0	1
Direction	0°	-45°	-90°
			
Number of floors	1	4	7
Surface absorptance	0.40	0.60	0.80
Indoor temperature	72°F	74°F	76°F
Occupancy	320 ft ² /person	220 ft ² /person	120 ft ² /person
Ventilation	0.06 cfm/ft ²	0.12 cfm/ft ²	0.18 cfm/ft ²
Opaque surface R	R = 5.65 °F·ft ² ·hr/Btu	R = 6.89 °F·ft ² ·hr/Btu	R = 8.13 °F·ft ² ·hr/Btu
Window ratio	20%	35%	50%

Table 4-5. The DSD factor assignment for each simulation run.

Continuous Factors									Categorical Factors		
Simulation Run	Direction	Floors	Surface absorptance	T_{in}	Occupancy	Ventilation	Wall R	Window Ratio	Aspect Ratio	Glass Types	Cold deck T
1	0	1	1	1	1	1	1	1	2	2	2
2	0	-1	-1	-1	-1	-1	-1	-1	1	1	1
3	1	0	-1	1	-1	-1	-1	1	2	2	1
4	-1	0	1	-1	1	1	1	-1	1	1	2
5	1	1	0	-1	1	-1	-1	-1	2	2	2
6	-1	-1	0	1	-1	1	1	1	1	1	1

Table 4-5 Continued.

Continuous Factors								Categorical Factors			
Simulation Run	Direction	Floors	Surface absorptance	T_{in}	Occupancy	Ventilation	Wall R	Window Ratio	Aspect Ratio	Glass Types	Cold deck T
7	1	-1	1	0	-1	1	-1	-1	1	2	2
8	-1	1	-1	0	1	-1	1	1	2	1	1
9	1	1	-1	1	0	-1	1	-1	1	1	2
10	-1	-1	1	-1	0	1	-1	1	2	2	1
11	1	1	1	-1	1	0	-1	1	1	1	1
12	-1	-1	-1	1	-1	0	1	-1	2	2	2
13	1	1	1	1	-1	1	0	-1	2	1	1
14	-1	-1	-1	-1	1	-1	0	1	1	2	2
15	1	-1	1	1	1	-1	1	0	1	2	1
16	-1	1	-1	-1	-1	1	-1	0	2	1	2
17	1	-1	-1	1	1	1	-1	1	2	1	2
18	-1	1	1	-1	-1	-1	1	-1	1	2	1
19	1	-1	-1	-1	1	1	1	-1	2	2	1
20	-1	1	1	1	-1	-1	-1	1	1	1	2
21	1	1	-1	-1	-1	1	1	1	1	2	2
22	-1	-1	1	1	1	-1	-1	-1	2	1	1
23	1	-1	1	-1	-1	-1	1	1	2	1	2
24	-1	1	-1	1	1	1	-1	-1	1	2	1
25	0	0	0	0	0	0	0	0	1	1	1
26	0	0	0	0	0	0	0	0	2	2	2

4.2.3 Building Energy Model

A building energy model was developed using eQuest for the parametric runs with the following conditions. The inputs to this simulation model were modified, as in Table 4-5, to generate datasets for regression analysis.

- Constant indoor temperature

- Four perimeter zones (15-ft width), one core conditioned zone, and one plenum for each floor
- Single-duct VAV systems (one set for each floor)
- Occupancy, lighting and equipment with the same daily schedules throughout the year
- Constant outside airflow
- Four-inch heavy weight concrete walls and roof
- U-values varied by adjusting the insulation resistance
- Same floor area (31,258 ft²) regardless of the building dimensions
- Equal window ratios on all the four vertical surfaces
- Weather data from the TMY3 College Station
- Ground floor assumed adiabatic
- Shading factor = 0.5
- Lighting load = 1.11 W/ft², and plug load = 0.43 W/ft²

4.2.4 Model Selection

All possible regression models involving the main effects, two-factor interactions, and pure-quadratic terms were fit, and the optimal model from a group of parametric models was selected using the Bayesian Information Criterion (BIC) (Schwarz 1978). The model to be chosen is the one that minimizes

$$\text{BIC} = -2 \log(\text{Maximum Likelihood}) + k \log(n) \quad (4.1)$$

where k is the number of parameters, and n is the number of observations. This criterion penalizes the introduction of new parameters, and it overcomes overfitting problems.

The model fit and selection process using the stepwise regression technique was performed using JMP® (SAS 2014).

4.2.5 Analysis

The variabilities of the T_{ref} value and factor sensitivity on the model are analyzed using the optimal model based on the BIC. To estimate the possible ranges for the given factor levels, a Monte Carlo Simulation (MCS), which is a random sampling method, is employed with the assumption that all the factors have normal distributions. The MCS is an approximate inference method based on random sampling. For each input, the MCS procedure generates a random sample from a given probability distribution. Then, random samples for all inputs are used to produce a single output using the known functional relationship between the inputs and output. This process is repeated a large number of times, and the mean and standard deviation of these output values are respective estimates of the output quantity and the associated standard uncertainty.

Each of the continuous factors was assumed to have a normal distribution, with a mean at the center value of the factor level and a standard deviation of $a / \sqrt{3}$, where a is the half-width between the upper and lower limits. This is a common conversion of a rectangular distribution to a normal distribution for estimating measurement uncertainty (JCGM 2010). By using this conversion, the means (μ 's) and standard deviations (σ 's) of normal distributions for the continuous factors were defined, as in Table 4-6. For the categorical factors, the probability of each level was set to 0.5. A normal distribution is fitted to the results from 5,000 trials to estimate the approximate distribution of the T_{ref} parameter.

Table 4-6. Normal distributions assumed for continuous factors.
 μ = mean, σ = standard deviation.

Factor	μ	σ
Direction	-45°	25.9°
Floors	4	1.7
Absorptance	0.6	0.11
T_{in} ($^\circ\text{F}$)	74	1.1
Occupancy ($\text{ft}^2/\text{person}$)	220	57.7
Ventilation (cfm/ft^2)	0.12	0.034
Wall Roof R ($^\circ\text{F}\cdot\text{ft}^2\cdot\text{hr}/\text{Btu}$)	6.89	0.715
Window Ratio	0.35	0.086

To assess the effect of each factor on the response, the sensitivity indices that measure the importance of factors in a model were estimated using JMP[®] (SAS 2014). The indices include the main effect and the total effect, and they vary from 0 to 1, with 0 meaning no effect. If the variation in the factor causes a high variability in the response, then that factor is considered to be important relative to the model, and the index values increase. A brief summary of the background is presented next, and the statistical details can be found in SAS (2014) and Saltelli (2002). Given a mathematical model

$$y = f(x_1, x_2, \dots, x_k) , \quad (4.2)$$

the expected value of y , namely $E(y)$, is defined by integrating y with respect to the joint distribution of the input factors—the x_i 's—and the variance of y , namely $Var(y)$, is

defined by integrating $(y - E(y))^2$ with respect to the joint distribution of the x_i 's. The impact of the main effect x_j on y can be described by $Var(E(y | x_j))$, and the ratio

$$\frac{Var(E(y | x_i))}{Var(y)} \quad (4.3)$$

provides a measure of the sensitivity of y to the factor x_j . The total effect represents the total contribution to the variance of y from all terms that involve x_j , including the interactions. For example, if there are only two factors, x_1 and x_2 , then the total effect of the importance index for x_1 is an estimate of

$$\frac{Var(E(y | x_1)) + Var(E(y | x_1, x_2))}{Var(y)} \quad (4.4)$$

where (x_1, x_2) is an interaction term of x_1 and x_2 . The computation of variances for the indices use the Monte Carlo procedure.












4.3 Results and Discussion

The estimates of the final T_{ref} model can be found in Table 4-7, and the main effect and total effect indices for this final model can be found in Table 4-8. According to the total and main effect indices, ventilation, occupancy, and indoor temperature had the largest effects on T_{ref} , followed by the window glass type, the cold deck temperature reset, and the surface absorptance. The effects from building construction factors such as the direction, the window to wall ratio, the thermal resistance of walls, the number of floors, the aspect ratio of the floor shape, and the surface absorptance ranked low. According to the result from the MCS of this T_{ref} model, the mean of the T_{ref} value was estimated as $63.4^\circ\text{F} \pm 2.4^\circ\text{F}$ when the factors were varied across the distributions in Table 4-6. The 2σ range of the T_{ref} value based on this estimation is $58.6\text{--}68.1^\circ\text{F}$.

Table 4-7. The final T_{ref} model and estimates selected based on the BIC. ‘×’ means interactions of two factors.

Term	Estimate	Standard Error	P-value
Intercept	63.4	0.058	< .0001
Direction	0.105	0.063	0.1420
Floors	−0.414	0.063	0.0003
Absorptance	−0.759	0.063	< .0001
Indoor T	1.73	0.063	< .0001
Occupancy	1.98	0.063	< .0001
Ventilation	2.06	0.063	< .0001
Wall R	−0.275	0.063	0.0035
Window ratio	−0.262	0.063	0.0045
Aspect ratio	−0.127	0.059	0.0704
Glass type	−0.860	0.059	< .0001
Cold deck reset	−0.397	0.059	0.0003
Floors×Indoor T	0.340	0.077	0.0033
Floors×Occupancy	−0.144	0.108	0.2261
Indoor T×Aspect ratio	−0.444	0.102	0.0034
Indoor T×Glass type	−0.550	0.083	0.0035
Occupancy×Ventilation	−0.286	0.115	0.0427
Wall R×Window ratio	−0.162	0.088	0.1094
Glass type×Cold deck reset	−0.887	0.126	0.0002

Table 4-8. The main effect and total effect indices for the final T_{ref} model. The bar chart illustrates the size of the total effect index.

Column	Main Effect	Total Effect	
Ventilation	0.170	0.367	
Occupancy	0.157	0.343	
Indoor T	0.130	0.310	
Glass type	0.075	0.217	
Cold deck reset	0.055	0.161	
Absorptance	0.039	0.086	
Aspect ratio	0.028	0.066	
Floors	0.028	0.063	
Wall R	0.015	0.032	
Window ratio	0.015	0.031	
Direction	0.004	0.009	

The sensitivity study demonstrated that the ventilation rate and occupancy level are the two largest effects, followed by the indoor temperature, which has not been known before. This can be explained by analyzing the functional expression of E_{BL} . If $T_{oa} = T_{in}$ in Equation (3.2), with the parameters described in Table 3-1, then E_{BL} becomes independent of T_{oa} . The E_{BL} value when $T_{oa} = T_{in}$ can be written as follows:

$$E_{BL} \big|_{T_{oa}=T_{in}} = -m_v h_v W_{oa}^+ - Q_{occ} - Q_{sol}. \quad (4.5)$$

The E_{BL} value when $T_{oa} = T_{in}$ can be viewed as a stationary point if drawn in an E_{BL} vs. T_{oa} plot, as depicted in Figure 4-2. Here, ‘stationary’ means that the point does not move, regardless of changes in the T_{oa} slope. The E_{BL} value at the stationary point $E_{BL}|_{T_{oa}=T_{in}}$ is always negative, meaning that the point stays below the T_{oa} axis. Furthermore, T_{ref} represents the T_{oa} value when $E_{BL} = 0$, and in the E_{BL} vs. T_{oa} graph, T_{ref} can be viewed as the point where the E_{BL} line crosses the T_{oa} axis. The graphs in Figure 4-2 demonstrate that T_{ref} is more sensitive to the slope as the stationary point moves farther below the T_{oa} axis. According to Equation (4.5), the stationary point will go down farther as the outdoor latent load increases, the occupant load increases, and/or the solar load increases. This explains why the ventilation and occupancy can be the factors that have the highest effects on T_{ref} . According to this observation, highly ventilated buildings should have a higher T_{ref} that is close to the indoor temperature.

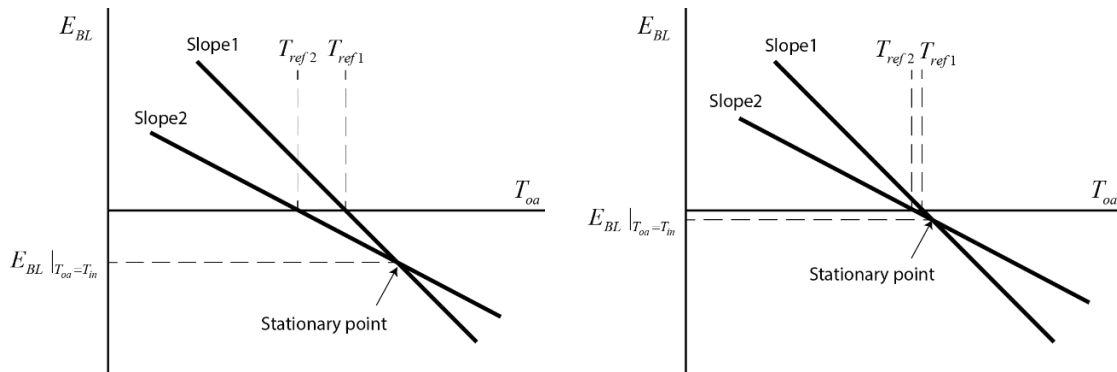


Figure 4-2. Impact of the level E_{BL} vs. T_{oa} stationary point on T_{ref} .

The rule of thumb for the T_{ref} values that have been used in the E_{BL} analysis is approximately 70°F. The distribution of the T_{ref} estimates using the MLR(T,W) model

for the energy use dataset in Chapter 2, and the data collected from 65 buildings on the Texas A&M University campus, was $66.7^{\circ}\text{F} \pm 3.7^{\circ}\text{F}$. The T_{ref} distribution based on this sensitivity and uncertainty analysis was $63.4^{\circ}\text{F} \pm 2.4^{\circ}\text{F}$, which was lower than the distribution of the estimates based on actual data. The actual building parameters are unknown, and the cause of this discrepancy cannot be clearly identified. However, some differences between the analysis conditions and the actual buildings can be pointed out as possible reasons. Some buildings have higher space temperature set points at night, and this can increase the daily average indoor temperature above that assumed in the factor range. The electricity use in the actual buildings may include some cooling energy use for window units and DX units in some buildings. Moreover, laboratory buildings have much higher ventilation rates than the high limit of the range in the factor analysis, and the blinds in actual buildings may be more on the closed side than the simulation. All of these elements can increase the discrepancy between the simulated and the actual T_{ref} values.

The limitation of this study is that all the factors are assumed to be independent of each other. For most of the factors, this assumption can hold; however, the occupancy level and ventilation rate should be correlated for actual buildings, and some corrections may be necessary. For example, in the T_{ref} model derived from the analysis, the lowest T_{ref} level occurs when the occupancy density is at the highest level and the ventilation rate is at the lowest level. This is not likely to be true in actual buildings because the ventilation rate usually increases as the occupancy density increases. It is possible that

the low limit of the T_{ref} range is underestimated because of this misspecification of the model.

The sensitivity and uncertainty analysis on the T_{ref} in this chapter presented the method to estimate the T_{ref} range when multiple influential factors vary at the same time. The study confirmed that there is a range that T_{ref} can take for the group of campus buildings, and the ventilation rate and occupancy level are found to be as influential as the indoor temperature. Using the T_{ref} model derived in the factor analysis, one can predict the degree to which T_{ref} can vary for different factor levels. This is useful in developing recommended T_{ref} ranges for different types of buildings.

5. CONTINUOUS BUILDING ENERGY DATA MONITORING USING RECURSIVE LEAST SQUARES FILTER AND CUSUM CHANGE DETECTION: APPLICATION TO ENERGY BALANCE LOAD DATA³

This chapter investigates a data-driven analysis method to detect abnormal energy data using the recursive least squares (RLS) filter and the cumulative sum (CUSUM) test. This model-based method compares the value predicted by the RLS filter and the actual value, and the CUSUM test sounds an alarm if the difference exceeds the prescribed threshold. In the present work, the method is applied to the whole building E_{BL} analysis using the outside air temperature and latent load variables on a daily basis. The ratios of the RMSE of the RLS filters to the RMSE of the regression solutions for 15 sample buildings during a one-year period range from 0.69 to 0.97. In the two case studies, the temperature drift of a chilled water meter was detected on the fourth day, and the disabled occupied/unoccupied HVAC schedule was detected on the seventh day after the problems appeared. Updating reference models to account for the dynamic use and operations of buildings has been a challenge in the implementation of the existing model-based fault-detection methods. The proposed method can track time-varying parameters automatically, and it requires less effort to maintain the prediction performance of the reference models.

³ Reprinted with permission from ASHRAE Transactions, Vol.121, Part1, Hiroko Masuda and David E. Claridge, Continuous Building Energy Data Monitoring Using Recursive Least Squares Filter and CUSUM Change Detection: Application to Energy Balance Load Data, pp.361-373. Copy right 2015 American Society of Heating, Refrigerating and Air-Conditioning Engineers, Inc.

5.1 Introduction

The continuous monitoring of building energy use offers useful feedback to the building owner and operators to achieve persistent building efficiency. Many software tools are available on the market to visualize and analyze energy use data, collected from acquisition systems or building automation systems, to support effective energy monitoring (Granderson et al. 2009; Ulickey et al. 2010). One important form of feedback that is obtainable from these energy analysis tools is the detection of abnormal increases and decreases in the energy use from various causes, including metering errors, equipment failure, energy efficiency improvements, and planned or unplanned changes in the HVAC controls and schedules. The timely detection of such changes can be used for effective energy management.

The detection of abnormal system changes is referred to as change detection, and if the detection involves malfunction(s) and/or performance degradation, this aspect of change detection is referred to as fault detection (Patton et al. 2000). A change-detection algorithm commonly consists of a residual generator and a change detector (Gustafsson 2001). The residual generator takes the difference between the model prediction and the actual value, and the change detector sounds an alarm if the residual exceeds the predefined threshold, which indicates an abnormal change in the input signal. This change-detection scheme is also known as model-based fault detection (Isermann 1997) or analytical redundancy (Clark et al. 1975; Chow and Willsky 1984).

Several authors have studied change-detection methods for whole-building energy use in the context of fault detection and diagnosis for HVAC systems. Among those, there are two types of reference models: calibrated simulation models and statistically estimated, data-driven models. Maile et al. (2012) and O'Neill et al. (2014) utilized the measured data from building automation and control systems (BACS) as input to a calibrated EnergyPlus energy simulation, so the model can emulate actual energy use behavior in real time once it is automated. These tools can provide useful information for improvements and problem isolation; however, the modeling and calibration of detailed simulations require high-level skills, along with comprehensive knowledge of buildings, to achieve accurate results. Bynum et al. (2012) aimed at low-cost fault detection that can be implemented in existing buildings using limited measurements and simpler simulation at the trade-off of diagnostics granularity. Several fault-detection tools using data-driven models have been reported, including those of Haberl and Claridge (1987) and Dodier and Kreider (1999). Haberl and Claridge (1987) used linear regression models as functions of weather and other building-specific variables, while Dodier and Kreider (1999) employed artificial neural networks, which are a class of non-linear regression models that use weather and calendar variables. Data-driven models do not require calibration effort and detailed knowledge of the buildings. For this reason, regression models using weather and other limited variables are featured in many of the advanced energy information systems (EISs) on the market as a means of prediction and energy tracking, so the users can detect abnormal energy use (Granderson et al. 2009; Friedman et al. 2011; Kramer et al. 2013).

Updating the reference models is a common challenge in implementing these model-based change-detection methods for energy data. Buildings occasionally have renovations of and updates to the systems, a change of tenants and space functions, or a change in the HVAC controls, and these cause changes in the energy use. After these changes, the current model becomes obsolete, and one would need a new reference model to be able to detect abnormal energy patterns. Then, re-calibration is required for a simulation model (Bynum et al. 2012), and the collection of data to estimate a new model is needed for a data-driven model (Dodier and Kreider 1999) in order to predict normal energy use in the new state. In practice, these changes can occur gradually and/or frequently, which makes it difficult to maintain an effective change-detection scheme.

This chapter proposes a change-detection method for energy data, in which a time-varying reference model is automatically updated. The exponentially weighted recursive least squares (RLS) filter is used for the model estimation and residual generation, and the cumulative sum (CUSUM) test is used as a change detector. The RLS filter can track slow changes in the model parameters by discounting older information, and it works as a low-pass filter. When abrupt changes occur in the model parameters, the slowly adapting filter increases the residuals, which will be detected by the CUSUM test.

Hilliard and Jamieson (2013) proposed the use of recursive estimates charts as a supplement to CUSUM charts for energy monitoring and tracking applications. The recursive estimates with exponential memory loss, which is equivalent to the adaptive RLS filter, was applied to a linear, natural gas consumption model for a healthcare

facility as a function of heating degree-days and weekday variables. Over the course of three years, three changes were detected that are explained by an increase in ventilation air, a possible inefficient boiler operation, and a consumption increase from a new addition to the building. The main methodological difference between the present work and that of Hilliard and Jamieson (2013) is the change-detection method. Our method uses a sequential CUSUM test to detect significant changes in the residual mean of the RLS filter output, and CUSUM statistics will be reset after the detection of a significant change. On the other hand, Hilliard and Jamieson (2013) use the breakpoints estimation algorithm of Zeileis (2003), based on a class of the generalized fluctuation test, to retrospectively estimate the time of parameter changes in a certain span of data.

In the present work, the method is applied to energy data analysis using the EBL (Shao and Claridge 2006). The EBL data for the buildings on the Texas A&M University campus are used to test the proposed method. Two types of detected changes, namely a temperature reading drift in the chilled water energy meter and schedule changes in the HVAC operations, are presented.

5.2 Methodology

5.2.1 Energy Balance Load Model

In Chapter 2, it was demonstrated that the MLR model with the outside air temperature T_{oa} and the humidity variable W_{oa}^+ have a generally good fit. This chapter uses the MLR(T,W) model from Chapter 2:

$$E_{BL} = \beta_0 + \beta_T T_{oa} + \beta_W W_{oa}^+ + \varepsilon \quad (5.1)$$

where β_0 is the intercept, β_T and β_W are the parameters of T_{oa} and W_{oa}^+ respectively, and ε is an error term.

5.2.2 Overview of Data Analysis Method

Figure 5-1 illustrates the model-based change-detection scheme in the present work, based on classical fault-detection or change-detection schemes (Isermann 1997; Gertler 1991; Chow and Willsky 1984; Gustafsson 2001). In each instance, the RLS adaptive filter recursively updates the model estimates. The residuals are quantities that represent the inconsistency between the actual value and the model. If there is no change in the system, and if the model is correct, then the residuals are so-called white noise, that is, a sequence of independent stochastic variables with a zero mean and a known variance (Gustafsson 2001). The change detector sounds an alarm when a change is detected by a statistical test, based on the whiteness of the residuals. In the present work, the CUSUM test was used for the change detector, as described in Hinkley (1971), Benveniste and Basseville (1984), and Gustafsson (2001). The RLS filter and the CUSUM test are described in the following sections.

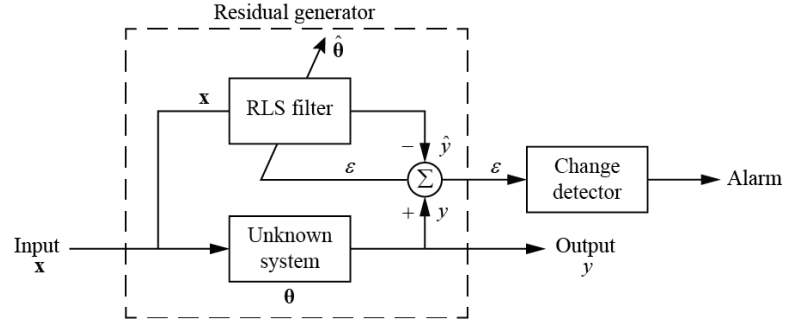


Figure 5-1. Change-detection scheme. \mathbf{x} = input, y = output, $\boldsymbol{\theta}$ = unknown parameters, $\hat{\boldsymbol{\theta}}$ = parameter estimates, and ε = residual.

5.2.3 The RLS Filter

We consider a time-varying regression model:

$$y(t) = \mathbf{x}^T(t)\boldsymbol{\theta}(t) + e(t) \quad (5.2)$$

where $y(t) \in R$ and $e(t) \in R$ are the observation and the noise at time t respectively;

$\mathbf{x}(t) \in R^N$ is a vector of explanatory variables, where N is the number of parameters; and

$\boldsymbol{\theta}(t) \in R^N$ is a vector of the unknown parameters to be estimated. The least squares solution aims to minimize a quadratic cost function, namely $J(\boldsymbol{\theta}) = e^2(t)$, to estimate the parameters $\boldsymbol{\theta}$. The exponentially weighted RLS, or simply RLS, approximates the loss function using the exponentially weighted sum of the residuals as:

$$J(\boldsymbol{\theta}) = \sum_{j=1}^t \lambda^{t-j} [y(j) - \hat{y}(j)]^2 = \sum_{j=1}^t \lambda^{t-j} [y(j) - \mathbf{x}^T(j)\boldsymbol{\theta}]^2 \quad (5.3)$$

where $0 < \lambda \leq 1$ is the forgetting factor, which offers a larger weight to recent data to track time-varying parameters. If $\lambda = 1$, then Equation (5.3) has infinite memory, and the criterion is equivalent to that of the linear regression solution. With a larger λ , the

tracking speed of the algorithm will be longer, and the estimation noise will be smaller. With a smaller λ value, the algorithm will have faster adaptation to system dynamics; however, the estimation noise will be larger. The derivation of the recursive estimation for this problem can be found in textbooks such as Ljung (1999), Haykin (2001), Gustaffson (2001), and Young (2011). The standard RLS algorithm updates the parameter vector $\hat{\boldsymbol{\theta}}(t) \in R^N$, the Kalman gain vector $\mathbf{g}(t) \in R^N$, and the covariance matrix $\mathbf{P}(t) \in R^{N \times N}$ for each time instance t , as in Equation (5.4).

$$\begin{aligned}\hat{\boldsymbol{\theta}}(t) &= \hat{\boldsymbol{\theta}}(t-1) + \mathbf{g}(t)[y(t) - \mathbf{x}^T(t)\hat{\boldsymbol{\theta}}(t-1)] \\ \mathbf{g}(t) &= \mathbf{P}(t)\mathbf{x}(t) = \frac{\mathbf{P}(t-1)\mathbf{x}(t)}{\lambda + \mathbf{x}^T(t)\mathbf{P}(t-1)\mathbf{x}(t)} \\ \mathbf{P}(t) &= \lambda^{-1}[\mathbf{P}(t-1) - \mathbf{g}(t)\mathbf{x}^T(t)\mathbf{P}(t-1)]\end{aligned}\tag{5.4}$$

A common practice is to use a zero vector for $\boldsymbol{\theta}(0)$ and to set $\mathbf{P}(0) = \alpha\mathbf{I}$, where \mathbf{I} is a unit matrix, and α is a large positive constant (say, 10^6 in general) (Young 2011). For the E_{BL} regression model in Equation(5.1), $y(t) = E_{BL}(t)$, and $\mathbf{x}^T = [1, T_{oa}(t), W_{oa}^+(t)]$.

Forgetting Factor. The forgetting factor can be interpreted by relating it to the number of samples that yield the same averaging effect on the exponential weights with the forgetting factor. This equivalent sample size can be derived by studying the sum of a geometric sequence, as in Brown (1963). The equivalent sample, size M of the forgetting factor λ for the exponentially weighted RLS is as follows (Gustafsson 2001):

$$M = (1 + \lambda) / (1 - \lambda) .\tag{5.5}$$

This value is also known as the average age of the data, and it has a close relation to the integrated moving average models (Box et al. 2008) and exponentially weighted smoothing (Brown 1963). In the field of control, it is customary to associate the

forgetting factor with the time constant of a low-pass filter. The time constant τ of an exponentially weighted filter is the sample size for which an impulse input decays to $e^{-1} \approx 36\%$ of its original size, and it can be approximated as follows (Mathews and Douglas 2001):

$$\tau = 1 / (1 - \lambda) . \quad (5.6)$$

This value is called the asymptotic sample length (ASL), and it indicates the number of important previous samples contributing to the estimation (Clarke 1985). The value is also known as the memory time constant (Ljung 1999), and it is used as a rule of thumb for the choice of the forgetting factor.

To evaluate the tracking performance of RLS with different values of λ , this study selected test data with a relatively stable E_{BL} pattern and applied the RLS filter. The data period is December 1, 2012–December 31, 2013 (13 months), and the 15 buildings in Table 5-1 were selected because the E_{BL} data for these buildings as a function of T_{oa} have stable patterns and no outliers. No outstanding issue was found for these buildings in the monthly data verification process. Buildings with different E_{BL} levels and patterns for weekdays and weekends were excluded because the model in Equation (5.1) cannot be applied to such data without modification.

The value of λ was varied from 0.7 to 1.0, with a 0.01 increment; furthermore, the RLS filter was estimated for each λ , and the λ having minimum RMSE was recorded for each building. The residual used in this process is the one-step-ahead prediction error, defined as:

$$\varepsilon(t) = y(t) - \mathbf{x}^T(t) \hat{\boldsymbol{\theta}}(t-1) . \quad (5.7)$$

The first month (December 1–December 31, 2012) was used for the learning period, and the RMSEs were estimated for the period January 1, 2013–December 31, 2013. For comparison, the RMSEs of the MLR models using the data during the period January 1, 2013–December 31, 2013 were computed.

The minimum values of the RMSE for the RLS estimates were compared with those of the MLR estimates in Table 5-1, and the RMSE versus λ is plotted for three selected buildings in Figure 5-2. Table 5-1 demonstrates that RLS produces a lower value of RMSE than MLR does for all 15 buildings, with the ratio of $RMSE_{RLS}/RMSE_{MLR}$ varying between 0.69 and 0.97. The lower ratio indicates that some parameter variations are captured by the RLS filter during the modeling period. The higher ratio suggests that the parameters may be stable during the modeling period, and there is no significant difference between the RLS and MLR prediction errors. For the 15 buildings, the range of the forgetting factors that yield a minimum RMSE is 0.83 to 0.98, and the median is 0.90. Based on this result, $\lambda = 0.90$ is used in the case studies. Figure 5-2 indicates that RLS generally provides better prediction than MLR, even with λ above 0.90. Figure 5-2 also demonstrates that RLS provides a lower value of RMSE for a rather broad range of forgetting factors, suggesting that the use of a forgetting factor for a larger group of buildings provides results that are superior to those involving MLR. According to Equations (5.5) and (5.6), the equivalent sample size is 19 days, and the time constant of the filter is 10 days when $\lambda = 0.90$.

Table 5-1. Forgetting factors having minimum RMSE (λ_{\min}) for 15 selected buildings. The minimum RMSE from RLS (RMSE_{RLS}) and RMSE from MLR (RMSE_{MLR}) are compared.

Building Number	λ_{\min}	RLS Minimum RMSE, Btu/(day·ft ²)	MLR RMSE, Btu/(day·ft ²)	RLS Minimum RMSE, kJ/(day·m ²)	MLR RMSE, kJ/(day·m ²)	Number of <i>EBL</i> Data during 2013	Ratio of RMSE _{RLS} /RMSE _{MLR}
1	0.94	14.6	16.6	166	189	365	0.88
2	0.87	18.0	22.6	204	256	365	0.80
3	0.92	66.2	95.3	751	1082	364	0.69
4	0.89	24.1	26.9	274	306	365	0.89
5	0.91	17.2	21.0	195	238	365	0.82
6	0.83	28.6	35.0	325	397	364	0.82
7	0.85	31.7	32.7	360	371	362	0.97
8	0.90	13.5	14.9	153	170	365	0.90
9	0.85	12.7	15.6	145	177	365	0.82
10	0.91	24.5	28.3	279	321	311	0.87
11	0.91	19.0	24.8	216	281	365	0.77
12	0.94	32.2	33.2	365	377	365	0.97
13	0.98	30.7	31.6	349	359	364	0.97
14	0.89	15.7	18.3	179	208	365	0.86
15	0.90	12.7	16.5	144	187	364	0.77

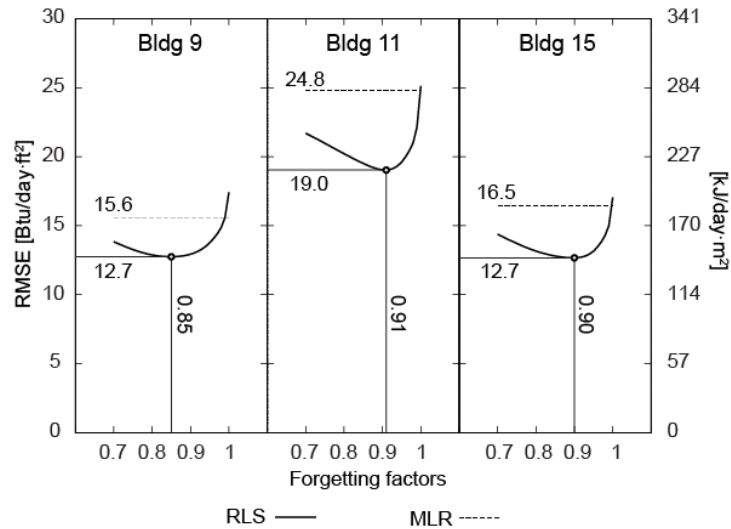


Figure 5-2. The RMSEs for RLS filter using varied λ for three buildings. The minimum RMSE values are indicated by points, and dashed lines represent RMSEs for MLR models.

Reset of Forgetting Factors. After the detection of changes, it is important that the filter will adapt to the new state quickly to continue monitoring data. One can control the forgetting factor to maintain the tracking capability of RLS. The mechanism for controlling this factor is to reduce the value of λ temporarily, so only new observations are used for estimation after parameter changes. There are many control methods for the forgetting factor (Leung and So 2005; Paleologu 2008); however, they usually involve weighting constants in addition to the forgetting factor. In the present work, this study uses the method suggested by Young (2011) to control λ at the beginning of the estimation process, and expand its use to the control of λ after a change is detected. The equation is a simple exponential decay process as a function of the constant forgetting factor defined for the filter and the sample size. The variable forgetting factor $\lambda(i)$ is

$$\lambda(i) = 1 - \frac{1 - \lambda}{1 - \lambda^{i+1}} \quad (5.8)$$

where $i = 1, 2, \dots$ is the time or sample size after a change is detected, and λ is the value that $\lambda(i)$ approaches as i becomes infinite. The value of i starts from 1 at the beginning of the estimation process.

Normalization of Residuals. The normalization of residuals allows one to use the same design of CUSUM detection for different buildings. The standardized residuals, which are the residuals scaled to the standard errors, are often used for the normalization of CUSUM control charts (Montgomery 2009). Brown (1963) described the derivation and examples of recursive variance estimation using exponential smoothing. This recursive estimation method is widely used in practice; for example, in on-line

monitoring for process variance (Wortham et al. 1974; MacGregor and Harris 1993) and the control of system identification for time-varying signals (Leung and So 2005; Paleologu 2008). Recursive estimation of the residual variance σ^2 using exponential smoothing is (Young 2011):

$$\hat{\sigma}^2(t) = \hat{\sigma}^2(t-1) + p(t)[\varepsilon^2(t) - \hat{\sigma}^2(t-1)] \quad (5.9)$$

where $p(t)$ is a variable weighting factor to discount the poor initial estimates but progressively utilize later data with equal weighting based on Equation (5.8). We expand the use of this weighting factor control to the reset of $\hat{\sigma}^2$ after a change is detected. By replacing t with i , which is the time after a change detection, $p(i)$ is defined as (Young 2011):

$$p(i) = \frac{1}{\lambda(i)} \left[p(i-1) - \frac{p^2(i-1)}{\lambda(i) + p(i-1)} \right] \quad (5.10)$$

and $p(i)$ approaches $(1-\lambda)$ as i increases. In the present work, $p(0)$ is set to 0.9. For $\lambda = 0.9$ and $p(0) = 0.9$, $|p(i) - (1-\lambda)| < 0.01$ is reached at $i = 34$.

By employing the residual variance estimator that is updated using Equation (5.9), the standardized residual s at time t is defined as:

$$s(t) = \varepsilon(t) / \hat{\sigma}(t-1). \quad (5.11)$$

The residual is normalized by the standard error one-step behind to have a better capability of detecting a sudden increase in the residual variance. This approach assumes that $\hat{\sigma}(t)$ values in adjacent steps are similar if there is no change.

5.2.4 The CUSUM Test

There are several algorithms used for sequential change detection in time series data, including the exponentially weighted moving average (EWMA) or the geometric moving average (GMA) algorithm introduced by Roberts (1959), the sequential probability ratio test (SPRT) based on Wald (1947), and the CUSUM algorithm proposed by Page (1954). These algorithms are widely used, and many modifications have been made to improve their performances for complicated and practical applications. In the present work, we use the CUSUM algorithm for its simplicity and ease of implementation.

The standardized, two-sided CUSUM statistics are defined as (Montgomery 2009):

$$\begin{aligned} C^+(i) &= \max[0, s(i) - k + C^+(i-1)] \\ C^-(i) &= \max[0, -k - s(i) + C^-(i-1)] \end{aligned} \quad (5.12)$$

where k is called the reference or the allowance, h is called the decision interval, and these are multiples of σ . The upper and lower test statistics, C^+ and C^- , sum the positive and negative values of the input $s(i)$. The alarm is raised when either C^+ or C^- exceeds h , which is a threshold of the CUSUM test, and C^+ and C^- are reset to zero. To prevent a false alarm, the allowance k is subtracted by the input at each time instance. The usual choice of k is 0.5, which is the appropriate choice for detecting a 1σ shift in the mean (Ryan 2011).

There are two important performance measures for a general statistical change detector. One is to quickly detect a change if it occurs, and the other is to raise a small

number of false alarms if there is no change. The average run length (ARL) function is defined as the mean time between alarms from the change detector as a function of the magnitude of the change (Gustafsson 2001), and it is used to evaluate both performance measures. It is desirable for the zero-state, in-control $ARL(0)$ (for the magnitude of change = 0) to be reasonably large, so that false alarms will not occur frequently. It is known that using $h = 4$ or $h = 5$ and $k = 0.5$ provides good ARL properties for a 1σ shift (Montgomery 2009; Ryan 2011), and we use $h = 5$ and $k = 0.5$ in the present work. The design of change detection using these values means that a residual smaller than 0.5σ at each time instance is ignored, and alarms will be raised if the CUSUM of the residuals exceeds 5σ in the upper or lower sides. According to the tabulated ARL values (Crosier 1986; Lucas and Crosier 2000), the two-sided ARLs in this condition are $ARL(0) = 465$, $ARL(1\sigma) = 10.40$, $ARL(2\sigma) = 4.01$, and $ARL(3\sigma) = 2.57$.

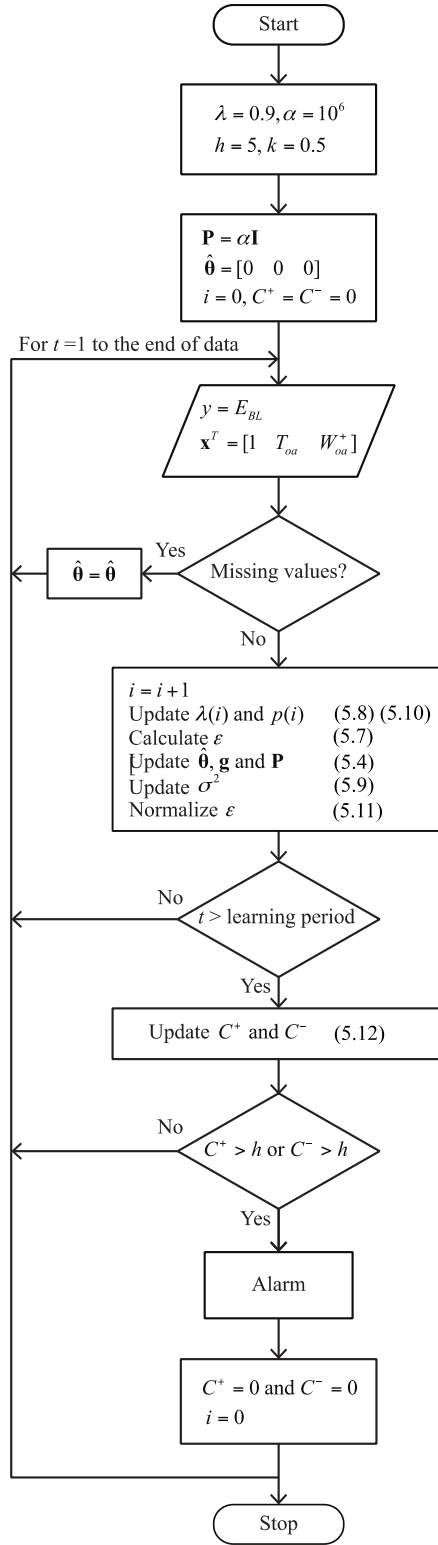


Figure 5-3. Flowchart for detecting changes in E_{BL} data

5.3 Results and Discussion

5.3.1 General Performance of RLS Filters and Change Detection for E_{BL} Data

Figure 5-3 summarizes the algorithm of the change-detection method. This algorithm was applied to the E_{BL} data from 149 campus buildings to understand the applicability and general performance. This section presents one example to explain the general performance found in the application of the change-detection method.

Figure 5-4 and Figure 5-5 are standard graphical outputs for one building. This particular set of plots presents the results for an office and laboratory building denoted as CHA (also as building 3 in Table 5-1) for the period of 2013. For RLS estimation, one additional month (December 2012) is used as the learning period; however, the CUSUM change detection is resumed on January 1, 2013. The E_{BL} has missing values for the first seven days, from December 1, 2012 through December 7, 2012. The primary result of the change detection is the plot of the CUSUM statistics with alarms in the bottom of Figure 5-4. The upper and lower CUSUM statistics C^+ and C^- are plotted in the positive and negative directions respectively to emphasize the increasing and decreasing trends. In Figure 5-4, the electric, cooling and heating energy use, the MLR residuals with approximate 2σ and 3σ prediction errors, and the RLS residuals are also presented. The parameter estimates and the standard error estimates of the RLS filter are plotted, along with the corresponding MLR estimates, in Figure 5-5.

Tracking Performance. In the given example, the time series of the MLR residuals has a positive bias during the first four months and a negative bias during the

last three months. Possible causes of these types of time-dependent variations in E_{BL} include changes in the indoor temperature, occupancy level, and control set points in HVAC systems. The RLS residuals do not have similar variations in the mean; instead, the parameters vary because of the adaptive algorithm. The standard error for the RLS filter is smaller than that for MLR, as can be seen in the bottom of Figure 5-5. The smaller prediction error increases the capability of detecting a change. In regression model-based energy tracking, such model bias can provide a wrong indication of an increase or decrease in energy use; and appropriate information regarding model uncertainties should be used, along with simple differences between model predictions and measured values. While the incorporation of auto-regression structures in a linear regression model can also reduce these variations in model residuals (Ruch, J. K. Kissock, et al. 1999; Liu et al. 2011), there is no guarantee that the same auto-regression structure exists in the energy data during the prediction period. The RLS filter cannot effectively track the E_{BL} using the simple model structure in Equation (5.1) if the building has quite different E_{BL} levels for weekdays and weekend days. For these buildings with apparent weekday/weekend patterns, the RLS residuals demonstrate almost identical patterns as the MLR residuals. Inclusion of day-type parameters (Hilliard and Jamieson 2013) or using separate filters can improve the predictions. One example of using separate filters for weekdays and weekends is provided later in this chapter.

Parameter Estimates. The parameter estimates of the MLR models for the E_{BL} variable have physical significance, such as ventilation rate and envelope thermal

performance (Masuda and Claridge 2012). The MLR parameter estimates approximately average the RLS parameter estimates for the same period, as illustrated in Figure 5-5, and the results suggest that the estimates from the RLS filter may also be physically interpreted. However, the fluctuation in the estimates, which is a characteristic of recursive estimates, makes it difficult to read the parameters. Increasing the value of λ will smooth the estimates to some extent; however, the adaptation gain in a linear adaptive filter is a compromise between noise attenuation and tracking ability (Gustafsson 2001), and the time delay obscures the interpretation of time-varying parameter estimates. Further study is needed for the capability and limitations of the parameter interpretation.

Convergence Speed. In the beginning of the RLS filter, the parameter estimates quickly approach the MLR estimates after the first several observations. The prediction error estimations using exponentially weighted smoothing also converge within the one-month learning period. This convergence performance is also applicable to the reset of λ after alarms, although the time required to converge depends on the length and the size of the level shift that triggered the alarm. Regardless of the types of level shifts, the parameter estimates converge within a period of several days to one month.

Cumulative Sum Change Detection. Most of the changes that are visually identifiable as outliers in the time series MLR residuals can be detected by the CUSUM change detection with the RLS filter. These changes in the E_{BL} are usually caused by metering problems, consumption drops due to temporary chilled water or heating water outage, and major changes in HVAC operations. Some of these cases are presented later

in this chapter. The proposed method also raise alarms for smaller level shifts. These changes might be related to occupancy levels, building use, and HVAC controls; however, such information is not available for the present study, and we are not able to verify the cause. The parameter estimates are disturbed, and the residuals increase after some alarms; however, the estimated standard errors also increase in such cases, which prevents repeating alarms from these disturbances. This, in turn, compromises the capability of detection right after alarms.

CUSUM Alarms for CHA. There are six alarms given by the CUSUM detection during January 1–December 31, 2013 for CHA. The first alarm is on March 14, 2013, after a steady increase in the C^+ statistic, which coincides with the decrease in the intercept and increase in the T_{oa} parameter. Many other buildings had alarms around the same time, and these might be related to the occupancy decrease during the spring break: March 11–March 15, 2013. The second alarm on March 29, 2013, is triggered by a sharp increase in the C^- statistic, which may have resulted from the decrease in the heating energy for two days during March 28–March 29, 2013. The heating hot water flow in the building was nearly zero for 22 hours during this period. The third alarm on August 25, 2013 is a result of the continuous increase in the C^- statistic starting from August 19, 2013. The fourth alarm on September 22, 2013 seems to be caused by a sharp decrease in the cooling energy use during September 21–September 22, 2013. The fifth alarm on November 6, 2013 is a result of a sudden increase in the C^- statistic after October 30, 2013. The 6th alarm on December 19, 2013 is from a sharp decrease in the C^- statistic, which is a result of the sharp drop in the heating hot water use during December 19–

December 20, 2013. The daily average outdoor temperature increased from 53.5°F (11.9°C) to 65.2°F (18.4°C) from December 18, 2013 to December 19, 2013; however, the decrease in the heating energy was larger than expected from the linear relationship with the temperature. In this manner, each alarm needs to be analyzed based on the available information to determine whether to dismiss it or not.

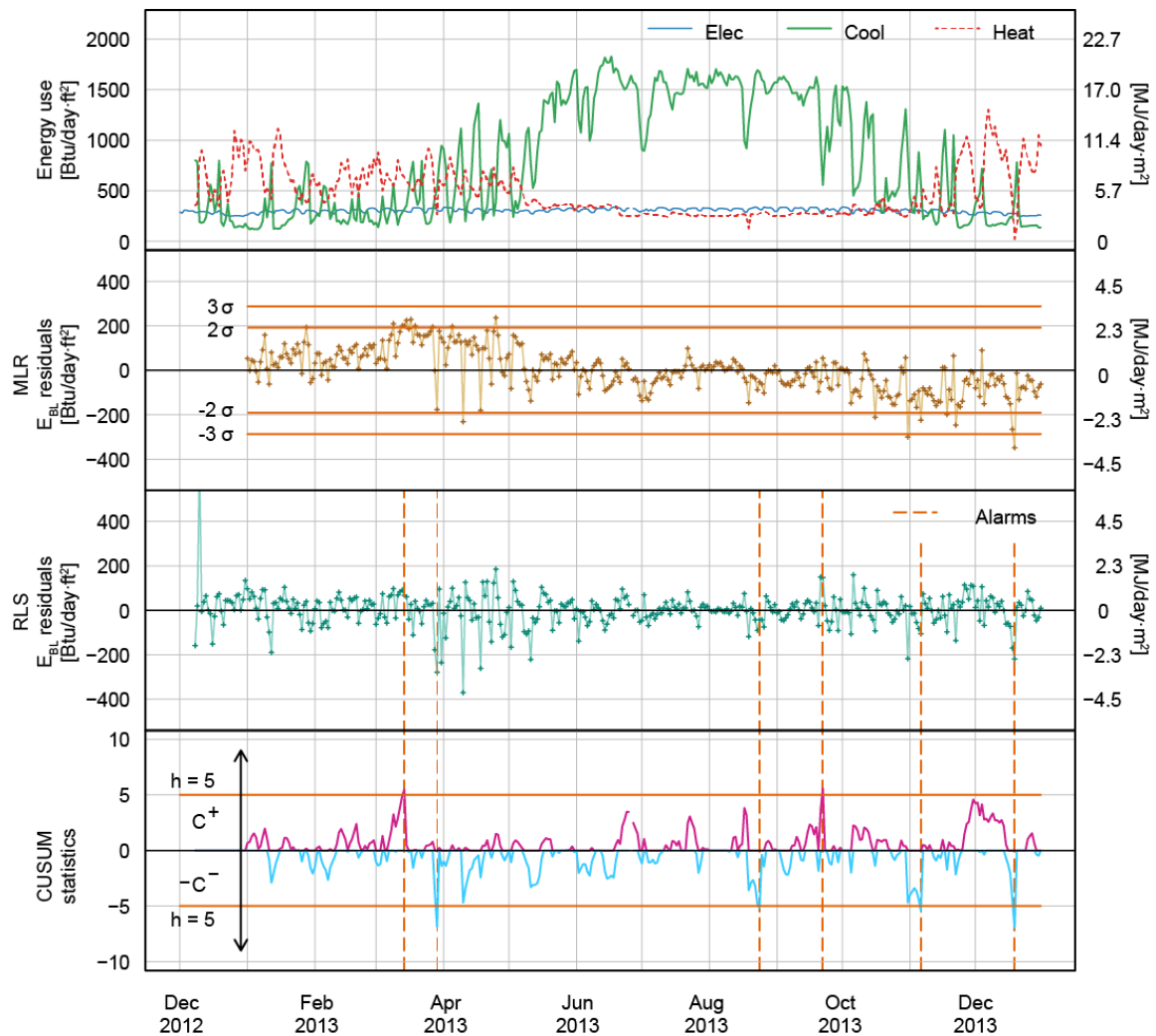


Figure 5-4. The result of RLS change detection ($\lambda = 0.9$, $h = 5$, $k = 0.5$) for CHA. The electricity (Elec), chilled water (Cool), and heating hot water (Heat) energy use and MLR model residuals are presented as well. Alarms are indicated as vertical lines.

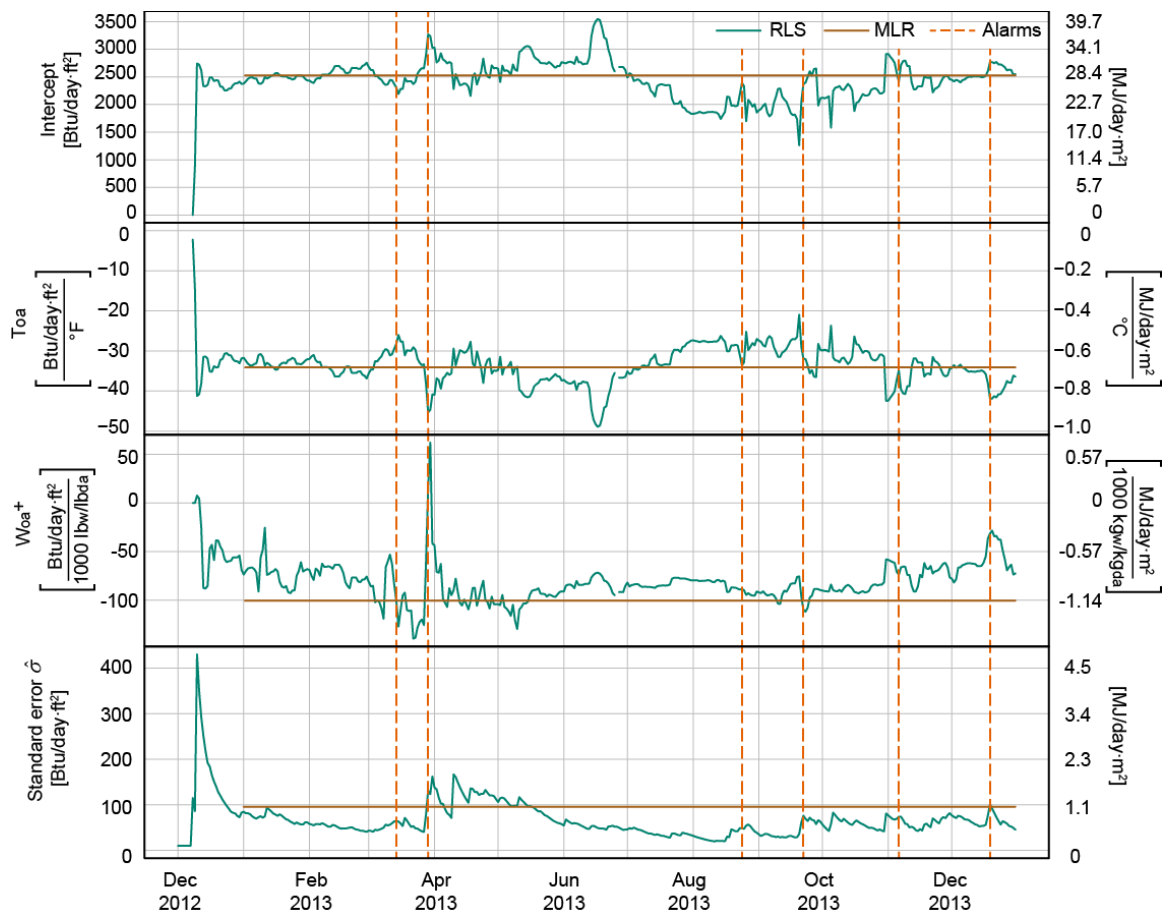


Figure 5-5. Parameter estimates and standard prediction error estimates of the RLS filter for CHA. The regression estimates for the MLR model are presented as well.

5.3.2 Detection of Metering Error from Temperature Measurement Drift

Temperature sensor drift is one of the frequently observed problems in thermal energy meters. The thermal energy use, such as chilled water and heating hot water, is calculated by multiplying the temperature differential between supply and return water and the flow rate, and the accuracy of the temperature sensors is important. Chilled water temperature differentials in particular for space cooling in commercial buildings can be small ($< 10^{\circ}\text{F}$ [5°C]), and the impact of temperature measurement bias is

significant. This case demonstrates how the proposed detection method detects the temperature drift problem in a chilled water meter.

The supply temperature for a library building, designated as MDL, started to decrease around June 30, 2013, as compared to those for some neighboring buildings. The chilled water for these buildings is supplied from the central plant, and the supply temperatures for these buildings are normally in the same range. Figure 5-6 compares the chilled water supply temperatures for MDL and one of the neighboring buildings. The largest temperature difference during this drift was approximately -2°F (-1.1°C). This type of temperature drift occasionally occurs in chilled water and heating hot water energy meters, and continuous effort is needed to detect and correct the drift.

In Figure 5-8, the level shift in the chilled water use can be visually seen in the energy use plotted as a function of the outside air temperature. However, the pattern during the last year has large variability, and it is difficult to differentiate the consumption level change due to a metering error from the variations in the normal activities and operations in the building. In the E_{BL} plot illustrated in Figure 5-9, the variability due to simultaneous cooling and heating in the building is removed, and the metering bias can be seen more clearly. In the energy balance plot, one can see the change in the pattern on July 25, 2013 and after.

The RLS filter and CUSUM change detection ($\lambda = 0.9$, $k = 0.5$, and $h = 5$) is applied to the E_{BL} data for this building. The filter estimation is resumed on December 1, 2012, and the CUSUM process is started on January 1, 2013. The increase of the C^- statistic resumes right after the temperature drift starts, and the first alarm is given on

July 4, 2013. The C^- statistic continues to decrease after the first alarm, which results in the second alarm on July 27, 2013. Then, the C^+ statistic increases suddenly on August 27, 2013, and the alarm is given on August 28, 2013. In this particular case, the supply temperature drift is detected about four days after the presumed start of the temperature drift. This is shorter than the detection time using either the E_{BL} versus T_{oa} plot or the chilled water consumption versus T_{oa} plot. The C^- statistic increases by $1.7 \sigma/\text{day}$, on average, from 0 on June 30, 2013 to 6.83 on July 4, 2013 and it exceeds $h = 5$ after three samples. This falls in the tabulated ARLs for the two-sided CUSUM process:

$ARL(2\sigma) = 4.01$ approximately.

When the start of the RLS filter is moved to May 1, 2013 (CUSUM starts on June 1, 2013), the alarm is still given on the same day, on July 4, 2013. When the start of the RLS filter is moved to June 1, 2013 (CUSUM starts on July 1, 2013), the alarm is delayed to July 10, 2013. This indicates that the method can detect abnormal changes with a short period of data. In this case, one to two months of data history is sufficient to detect the drift in a timely manner.

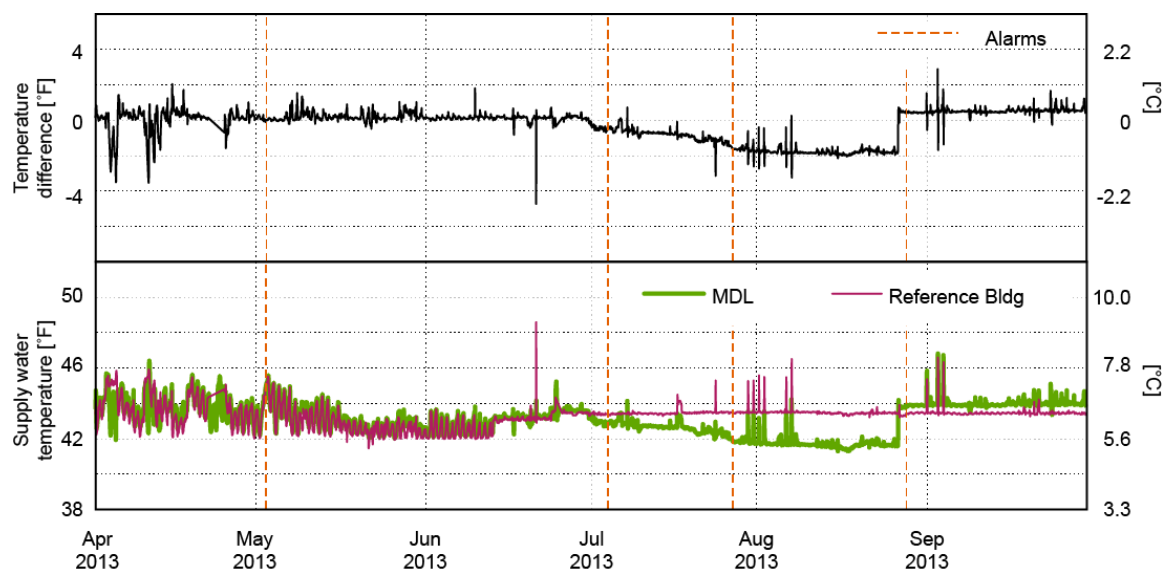


Figure 5-6. Hourly average supply water temperatures of the chilled water energy meters for MDL and a neighboring building from April 1 to September 30, 2013. The difference of the MDL meter from the adjacent building meter is plotted in the top portion, and the values for individual meters are plotted in the bottom portion of the figure. Alarms are indicated on 23:00 for each of the alarmed days.

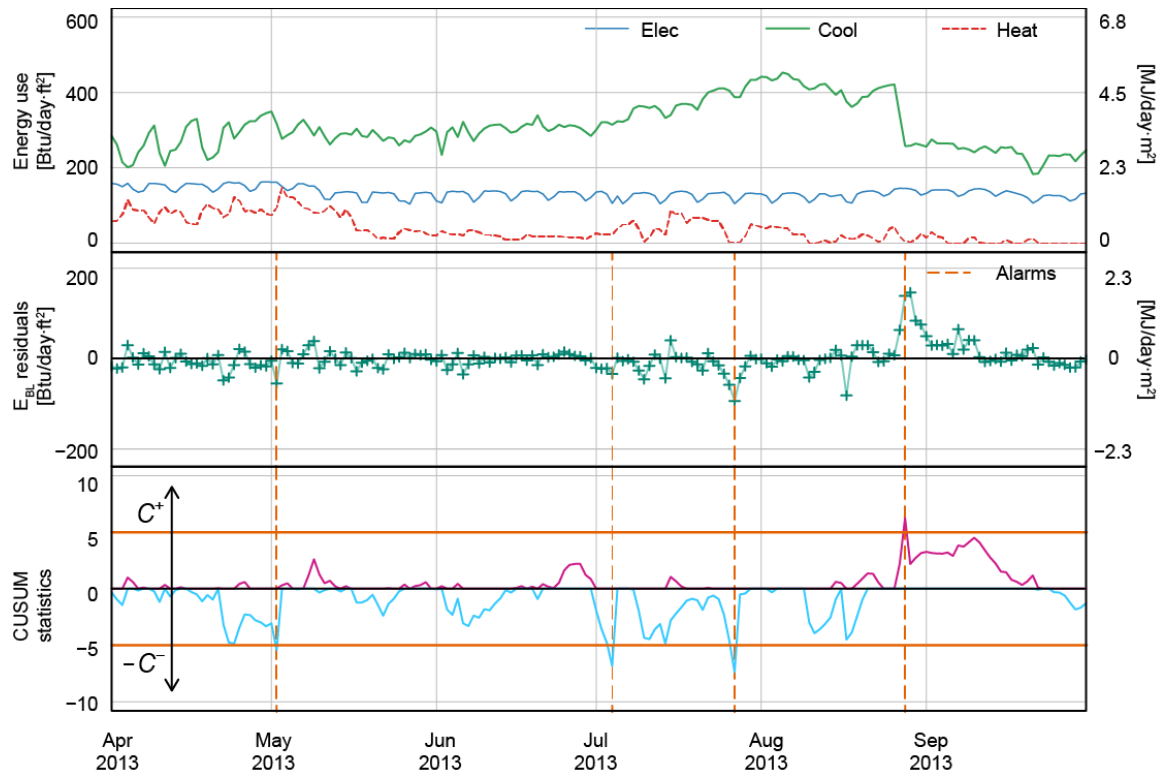


Figure 5-7. Energy use, RLS residuals, and CUSUM statistics for MDL during the period of April 1 to September 30, 2013 ($\lambda = 0.9$, $k = 0.5$, and $h = 5$).

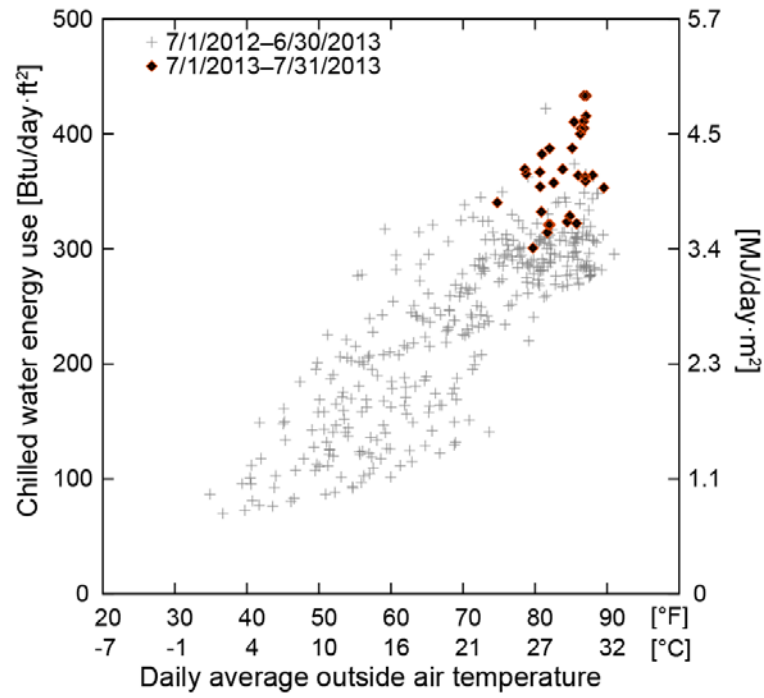


Figure 5-8. Chilled water energy use for MDL from July 1 to July 31, 2013 is compared to the values during the previous year.

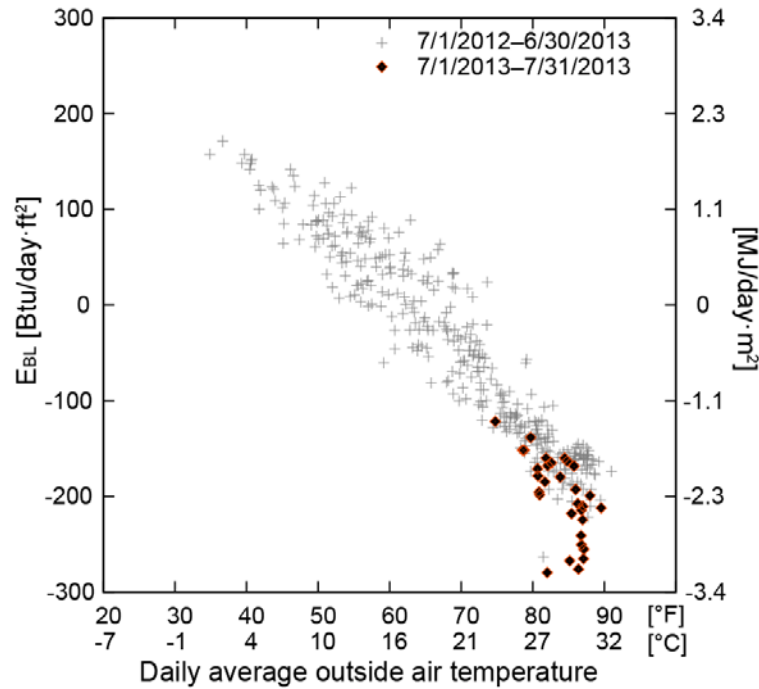


Figure 5-9. The E_{BL} for MDL from July 1 to July 31, 2013 is compared to the values during the previous year.

5.3.3 Detection of HVAC Schedule Change

For buildings with different weekday vs. weekend HVAC schedules that change E_{BL} levels significantly, a simple application of the linear model in Equation (5.1) is not appropriate. Figure 5-10 illustrates the E_{BL} versus T_{oa} plot for an office building, denoted as PVL, which has different patterns for weekdays and weekends. Although the HVAC scheduling information is not available, the E_{BL} and energy use patterns indicate the different HVAC operations for occupied and unoccupied hours. During the period of 7/8/2013-9/2/2013, the setback HVAC schedules for unoccupied hours did not appear to be implemented for some reason. The points during this period, highlighted in Figure

5-10, have the same pattern for weekdays and weekends, and they are in the same level as those for the weekdays during the other times of the year.

Two separate RLS filters and CUSUM detection processes are applied to weekday data and weekend data, and the results are presented in Figure 5-11. The data period is December 1, 2012 to December 31, 2013, with the first month as a learning period. The weekday residuals still have some periodic variations due to thermal transient effects from weekends to weekdays. For example, during the warm season, the E_{BL} values on Mondays tend to be lower due to extra cooling required to remove the heat added to the furniture and building mass by increased temperatures during the weekend. Some cold weeks during the heating season demonstrate the opposite changes in the E_{BL} on Mondays due to an extra heating load. A separate filter for Mondays decreases the residuals even more; however, only the result from two separate filters is presented.

The start of the schedule change is detected by the weekend process on July 14, 2013, which is the first Sunday after the schedule change began on July 8, 2013. The change-back of the schedule on September 3, 2013 is detected by the weekday process on September 4, 2013. The weekend process detects the change-back on September 14, 2013, which is the second weekend after the change. The weekday process does not detect the start of the schedule change well because the level does not change significantly, and the variance decreases. Meanwhile, it detects the change-back of the schedule fast because the variance increased.

This case study demonstrates that the separate processes can be used to analyze the data with different levels for different days of the week. The amount of data and the

sample frequency are different for each process, and additional consideration for the choice of λ and the design of CUSUM detection may be needed.

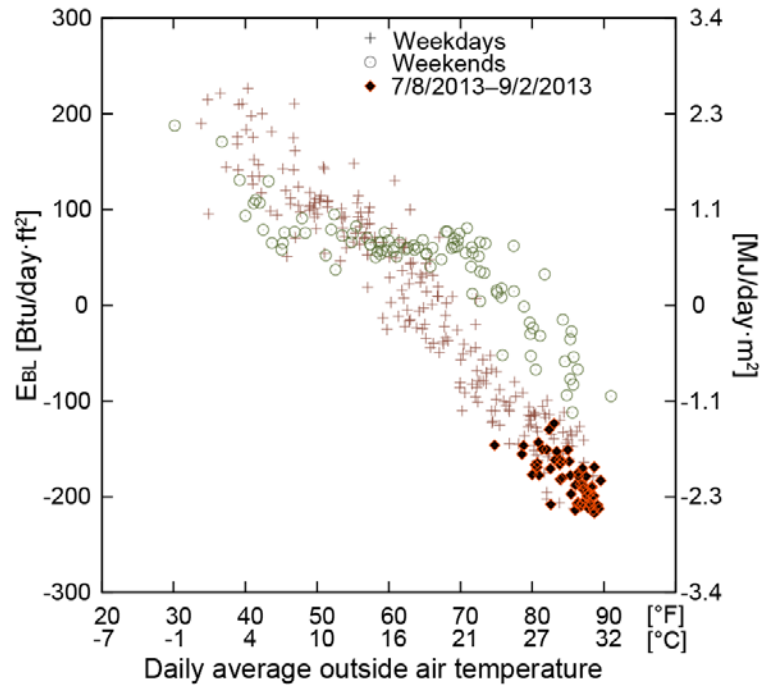


Figure 5-10. The E_{BL} for PVL from January 1 to December 31, 2013. Weekdays and weekends are plotted in different markers. The data between July 8 and September 2, 2013 are highlighted.

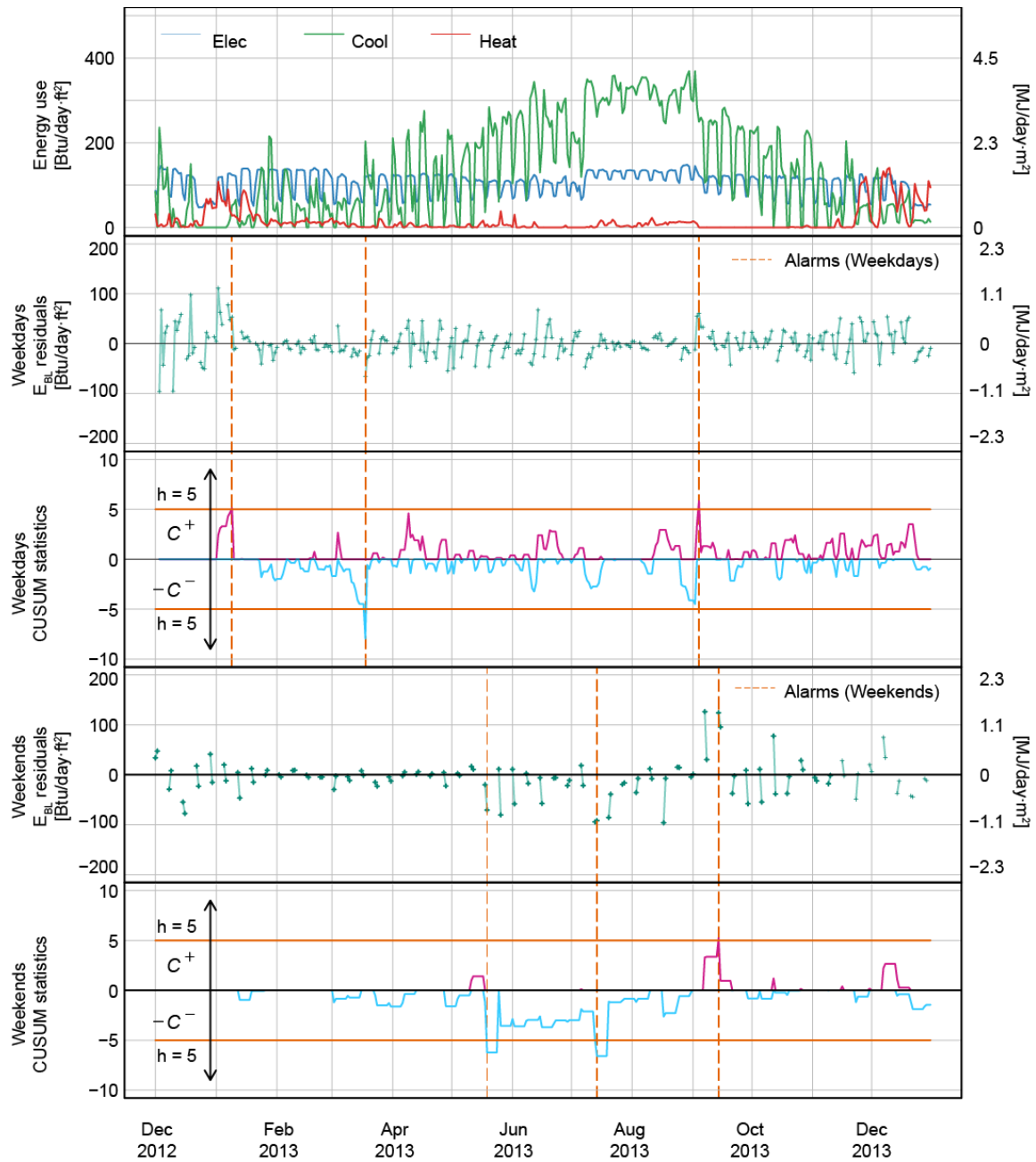


Figure 5-11. Energy use, RLS residuals, and CUSUM statistics of weekday and weekend processes for PVL during the period of December 1, 2012-December 31, 2013 ($\lambda = 0.9$, $k = 0.5$, and $h = 5$). The first month is the learning period.

5.3.4 *Limitations and Future Research*

Application to Energy Use Data. In the present work, the proposed change-detection method is applied only to the daily EBL data; however, it can also be applied to daily energy use data using appropriate linear models in the RLS filter. Several authors, including Fels et al. (1986), Rabl and Rialhe (1992), Reddy et al. (1995), Katipamula et al. (1994;1998), and Sonderegger (1998), proposed simplified linear regression model structures for cooling and heating energy use in a building, and these models could be used for the RLS filter estimation in the change-detection method. One challenging area for the future application to energy use data is the change point in these linear regression models. These cooling and heating models usually involve unknown change points that are unique to individual buildings, and the solutions iteratively search for them using the past data and regression estimates, assuming there is a fixed change point (Fels 1986; Schrock and Claridge 1989; Kissock et al. 2004). The recursive estimation does not search for unknown change points, and the proposed method may raise alarms for these change points unless they are already incorporated into the linear model for the RLS filter as known parameters.

Needs for Supporting Information. The primary output of the method is binary alarms to inform one of the occurrence of shifts exceeding the thresholds. The binary alarms can help energy managers and building operators to screen the data when large numbers of buildings are under continuous monitoring; however, they must decide whether to accept or dismiss the alarms based on other information. Effective

visualizations of energy data, model residuals, and CUSUM statistics are indispensable to support their decisions.

Verification of Estimated Models. While the proposed method detects changes, it does not consider the reasonableness of parameter estimates; therefore, a problem may not be detected if the parameters and variance are consistent from the beginning of the RLS filter estimation. To verify the estimated model, limit checks of the parameter estimates and assessment for the size of prediction errors could be integrated in the proposed method.

5.4 Chapter Summary

This chapter presented a method to monitor changes in energy data and its application to the E_{BL} data. The method used the adaptive RLS filter to automatically update the reference model to account for changes due to the dynamic use and operations of buildings. The standardized CUSUM test on the residuals was used for the change-detection process. The application of the method to the E_{BL} data can allow for timely alarms on some abnormal changes in the energy use. As examples, the detection of an energy measurement bias in the chilled water meter and HVAC schedule changes were presented.

The RLS filters can track variations in the parameters and decrease the prediction errors, compared to the regression solutions for the same period. For the E_{BL} data from 15 sample buildings during a one-year period, the ratios of the RMSE of the RLS filters to the RMSE of the regression solutions range from 0.69 to 0.97, depending on the stability of the parameters during the estimation period. The energy measurement bias in

the chilled water meter in MDL, caused by a drift of the temperature sensor reading, was detected on the fourth day after the temperature began to drift, whereas it can take 25 days until one can visually observe the pattern change in the E_{BL} versus T_{oa} plot. The temporary HVAC schedule change in PVL was detected using two separate filters and change-detection processes for weekdays and weekends. The start of the disabled occupied/unoccupied schedules was detected on the seventh day, and the change-back of the schedule was detected on the second day.

Updating the reference models to account for the dynamic use and operations of buildings is a challenge in the existing model-based fault-detection methods. The proposed method automatically updates the time-varying parameters, and it requires less effort to maintain the prediction performance of the reference model, compared to existing methods. This feature is suitable for a fully automated implementation as a part of energy tracking software tools.

6. CONCLUSIONS

The advancement of data collection and management technology has made energy use interval data available for a wide range of applications, including utility cost allocations for campus facilities, the assessment of energy efficiency projects, the tracking and benchmarking of building energy efficiency, and the calibration of building energy models. Despite increasing attention being paid to energy metering and data collection, stimulated by policies and energy cost, the quality control of collected energy data has not been widely discussed. Energy meters, especially thermal energy meters, often have errors that can cause misleading results in engineering and financial decisions, and it will take continuous effort to maintain the quality of metered energy use data. This dissertation attempts to contribute to the practice of managing the quality of measured energy use by developing a mechanism to automatically detect those errors with minimal prior information on buildings using a variable called the E_{BL} .

The objective of this dissertation is to develop a model-based detection method for anomalies in energy consumption data that are collected from a large number of buildings with minimal information using the E_{BL} variable. To pursue this goal, Chapter 2 developed data-driven E_{BL} models that can be applied to different types of buildings to describe the E_{BL} patterns numerically, Chapter 3 investigated the physical significance of the statistical estimates of the E_{BL} models, Chapter 4 investigated the sensitivity and uncertainty of T_{ref} , and Chapter 5 developed a new method for detecting energy use anomalies using the E_{BL} models and then demonstrated its application.

In Chapter 2, the structure of a steady-state E_{BL} model was derived as a linear combination of T_{oa} , W_{oa}^+ , E_{sol} , and E_{ele} , and four statistical models were proposed based on these variables. The applicability of these models was studied using the daily energy use data collected from 56 buildings on the Texas A&M University campus. The semi-partial correlations demonstrate that T_{oa} and W_{oa}^+ can explain most of the E_{BL} variance, and the MLR model with T_{oa} and W_{oa}^+ , designated as MLR(T, W), was found to be a simple yet widely-applicable E_{BL} model that could be fit to the E_{BL} for various types of buildings with readily available weather data. When applied to the data collected from 56 buildings, the mean of the model CV-RMSEs was 10.0%. The addition of the E_{sol} variable slightly decreased the mean of the model CV-RMSEs, from 10.0% to 9.7%. The inclusion of an autoregressive term decreased the mean of the model CV-RMSEs to 6.9%; however, the physical interpretation of the regression coefficients became difficult.

In Chapter 3, it is demonstrated that the parameter estimates for the MLR(T,W) model could be used for the approximate estimation of physical building parameters. Using the synthetic daily datasets from simulation models where the space temperature was maintained at its set point for more than 16 hours/day, the bias of the m_v estimates based on $\hat{\beta}_w$ was within 10%, and the bias of the T_{oa} slope estimates $\hat{\beta}_T$ was within 16%. The overall heat-loss coefficient estimates had higher biases in the range of 19% to 46% because the estimation involved two regression estimates: $\hat{\beta}_w$ and $\hat{\beta}_T$. This dissertation also compared the estimated and measured values of the air exchange rate for three dormitory buildings. Using E_{BL} daily data, the estimation biases were 2.5%, -12.0%, and

–27.8% respectively for the three buildings. The previous studies by Deng (1997) and Reddy et al. (1999) on the identification of physical parameters using synthetic data concluded that the method was accurate; however, based on the results from the actual data in this dissertation, one should be more cautious about the physical interpretation of parameter estimates. The accuracy of estimates depends on many factors, such as the quality of data, the constancy of building operation and controls, and the amount of data. Although the estimates are not always accurate enough to be used as measurements of physical values, the degree of physical significance demonstrated in this study indicates that it can be used to identify physically impossible parameter estimates. For example, abnormal estimates were indicated by the high value of T_{ref} , near 50°C (122°F) for the Hobby Hall monthly Q_B model, which was not a realistic value based on the physical significance of the parameter.

In the visual data-screening process using E_{BL} plots, the temperature at $E_{BL} = 0$ was used as an indicator of meter problems, based on the knowledge that the temperature must be close to the average indoor temperature; however, there was no quantitative study about the possible range of the value. In Chapter 2, this temperature was named the indoor reference temperature T_{ref} , and the estimation method was defined using regression coefficients. For the energy use data from 56 buildings, the mean and standard deviation of the T_{ref} estimates were 19.3°C and 2.0°C (66.7°F and 3.7°F respectively) using the MLR(T,W) model. In Chapter 4, the variability of the T_{ref} was estimated when influential factors varied within the range that represented a group of buildings on the Texas A&M University campus. The estimated 2σ range of the T_{ref}

value was between 58.6–68.1°F, with a mean of 63.4°F. The estimated mean was lower than the value for the actual dataset used in Chapter 2 by 3.3°F. Possible reasons for this discrepancy were discussed in Chapter 4. An important finding from this study is that one should expect that the T_{ref} estimate can vary about 10°F (5°C), depending on the type of building and HVAC controls for the given set of factor variations. The study also finds that the ventilation rate and occupancy level influence T_{ref} as much as the indoor temperature, indicating that highly ventilated laboratory buildings could have higher T_{ref} values, compared to dormitory buildings.

Chapter 5 developed a new model-based method to detect abnormal E_{BL} data. The new method used the RLS filter with a forgetting factor to estimate the E_{BL} prediction model as a function of T_{oa} and W_{oa}^+ , and anomalies were detected based on the size of the difference between the actual and predicted values using the CUSUM sequential test. With a forgetting factor λ , the newer data have more weight than the older data, providing adaptive estimation for slowly varying parameters. This approach has advantages in the continuous monitoring of a large number of buildings because the model will be automatically adjusted to up-to-date parameters, where model parameters vary over time due to changes in building use and operations. The ratio of RMSEs for RLS estimates with $\lambda = 0.9$ and MLR estimates ($RMSE_{RLS}/RMSE_{MLR}$) using annual data for 15 sample buildings ranged from 0.69 to 0.97, indicating that RLS estimates provide better prediction due to the adaptive estimation. A reset scheme for the estimates and forgetting factor after a change detection was incorporated so the model can be quickly adjusted for the new state. The new estimates can be stabilized within a month,

regardless of the types and magnitude of the detected changes. The reset scheme allows for continuous monitoring without the manual maintenance of prediction models. The method successfully detected a bias in the chilled water meter on the fourth day and a change in the HVAC system's time schedule on the seventh day after the problems began.

The new detection method demonstrates the ability to detect anomalies in energy use data in a timely manner. The model estimation is data-driven, and it requires only the T_{oa} and humidity ratio in addition to the energy use data. The model will be automatically updated to adapt to recent building use changes and operations. Therefore, the method provides an inexpensive way in which to continuously monitor energy use data for a large number of buildings.

6.1 Future Directions

The parameter estimates using RLS have physical significance, as in MLR estimates; however, their use was not studied in this dissertation. The interpretation of time-varying parameter estimates such as T_{ref} and the coefficients of T_{oa} and W_{oa}^+ could provide additional information to explain the detected changes. It was found that the RLS estimates fluctuate, and direct use of the values could be difficult. Some sort of smoothing may be necessary to utilize the RLS parameter estimates for further analysis.

The anomalies detected using this method can be not only meter errors but also the result of actual changes in the building use and operations. Future work should focus on incorporating the bottom-up approach, using the system-level information, such as

BAS data, to support the decision making regarding whether to accept or dismiss the detected changes.

6.2 Limitations

The E_{BL} is calculated from separately metered non-cooling electricity, cooling, and heating energy use. Therefore, if one energy source provides different types of end-uses—lighting, cooling, and/or heating—then the metered use has to be disaggregated to calculate the E_{BL} . For example, if the electricity is used for lighting and cooling, then the metered electricity use has to be disaggregated into lighting and cooling uses; if sub-metered electricity use is available, then one can accurately calculate the E_{BL} ; and if sub-meters are not available, this can still be done, with limited accuracy, using the known efficiency of the equipment.

The E_{BL} model in this dissertation was based on the assumption that the use of economizers and heat recovery was limited. The data used in this dissertation were from a hot and humid climate, and the impact of economizers and heat recovery on the whole building E_{BL} pattern was limited. Those interested in the impact of economizers and heat recovery on E_{BL} may wish to examine Shao and Claridge (2006) and Shao (2006).

The results of sensitivity and uncertainty studies on finding a possible range of T_{ref} in Chapter 4 depended on the selection of factor ranges. In this dissertation, the factor ranges were selected to describe a group of buildings on the Texas A&M University campus. The same approach to estimate the parameter ranges can be used for a different set of buildings; however, the different factor ranges will change the result.

The method is based on the assumption that the daily average of the indoor thermal condition is approximately constant. If the change is seasonal, then RLS can track the change; however, if the change is more frequent, such as daily or weekday/weekend, then the change is too rapid, and the method cannot effectively detect anomalies. If the changes follow a known schedule and produce clearly distinct levels, then this problem can be avoided by estimating a separate model for each level. A dummy variable may be added to the model, such as 0 for weekdays and 1 for weekends, to estimate different sets of parameters, and one could run two separate CUSUM processes. However, it takes longer for the RLS filter to learn the parameters if the data are separated. For example, the weekend model requires five weeks to obtain 10 days of data. This diminishes the advantage of the RLS and CUSUM detection methods.

REFERENCES

- Abushakra, B., Sreshthaputra, A., Haberl, J.S., and Claridge, D.E., 2000. *Compilation of Diversity Factors and Schedules for Energy and Cooling Load Calculations (RP-1093)*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers
- ASHRAE, 2009. *2009 ASHRAE Handbook: Fundamentals*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers.
- ASHRAE, 2017. *2017 ASHRAE handbook: Fundamentals*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers.
- ASHRAE, 2010a. *ANSI/ASHRAE/USGBC/IESNA Standard 189.1-2009: Standard for the Design of High-Performance Green Buildings Except Low-Rise Residential Buildings*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers.
- ASHRAE, 2014. *ASHRAE Guideline 13-2014: Specifying Direct Digital Control Systems*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers.
- ASHRAE, 2010b. *ASHRAE Guideline 2-2010: Engineering Analysis of Experimental Data*, Atlanta, GA: American Society of Heating, Refrigerating and Air-Conditioning Engineers.
- Baltazar, J.C., Sakurai, Y., Masuda, H., Feinauer, D., Liu, J., Ji, J., Claridge, D., Deng, S., and Bruner, H., 2007. Experiences on the Implementation of the “Energy Balance” Methodology as a Data Quality Control Tool: Application to the Building Energy Consumption of a Large University Campus. In *Proceedings of the 7th International Conference for Enhanced Building Operations*. San Francisco, CA: Energy Systems Laboratory. Available at: <http://hdl.handle.net/1969.1/6222>.
- Baltazar, J.C., Claridge, D.E., Ji, J., Masuda, H., and Deng, S., 2012. Use of First Law Energy Balance as a Screening Tool for Building Energy Data: Part II - Experiences of Its Implementation as a Data Quality Control Tool. *ASHRAE Transactions*, 118, pp.167–174.
- Baltazar, J.C., Haberl, J.S. and Sun, Y., 2011. *Solar Test Bench Database*, Energy Systems Laboratory, Texas A&M Engineering Experiment Station. Available at: <http://165.91.141.95:6785/>.
- Bauwens, G. et al., 2012. Reliability of co-heating measurements. In *First Building Simulation and Optimization Conference*. Loughborough, UK: IBPSA-England, pp. 417–424.

- Bauwens, G. and Roels, S., 2014. Co-heating test : A state-of-the-art. *Energy and Buildings*, 82, pp.163–172.
- Benveniste, A. and Basseville, M., 1984. Detection of abrupt changes in signals and dynamical systems : Some statistical aspects. In A. Bensoussan & J. L. Lions, eds. *Analysis and Optimization of Systems SE - 12*. Lecture Notes in Control and Information Sciences. Springer Berlin Heidelberg, pp. 143–155.
- Box, G.E.P., Jenkins, G.M. and Reinsel, G.C., 2008. *Time Series Analysis: Forecasting and Control* 4th ed., Hoboken, NJ: Wiley.
- Brown, R.G., 1963. *Smoothing, forecasting and prediction of discrete time series*, Englewood Cliffs, NJ: Prentice-Hall. Available at: <http://catalog.hathitrust.org/Record/001512022>.
- Butler, D. and Dengel, A., 2013. *Review of co-heating test methodologies*, November 2013, Milton Keynes, UK: NHBC Foundation/HIS BRE Press
- Bynum, J.D., Claridge, D.E. and Curtin, J.M., 2012. Development and testing of an Automated Building Commissioning Analysis Tool (ABCAT). *Energy and Buildings*, 55(0), pp.607–617.
- Capehart, B.L., Turner, W.C. and Kennedy, W.J., 2012. *Guide to Energy Management (7th Edition)*. Lilburn, GA: Fairmont Press, Inc.
- Casey, S., Krarti, M., Bianchi, M., and Roberts, D., 2010. Identifying Inefficient Single-Family Homes With Utility Bill Analysis Preprint. In *ASME 2010 4th International Conference on Energy Sustainability*. Phoenix, AZ.
- Childs, K.W., Courville, G.E. and Bales, E.L., 1983. *Thermal Mass Assessment: An Explanation of the Mechanism by Which Building Mass Influences Heating and Cooling Energy Requirements* (Report No. ORNL/CON-97). Oak Ridge, TN: Oak Ridge National Laboratory
- Choi, H., Yoon, B., Kim, C., and Choi, Y., 2011. Evaluation of flowmeters for heat metering. *Flow Measurement and Instrumentation*, 22(5), pp.475–481.
- Chow, E. and Willsky, A., 1984. Analytical redundancy and the design of robust failure detection systems. *Automatic Control, IEEE Transactions on*, 29(7).
- CIBSE, 2006. *TM41: 2006 Degree-days: theory and application*, London, UK: The Chartered Institution of Building Services Engineers
- Clark, R.N., Fosth, D.C. and Walton, V.M., 1975. Detecting Instrument Malfunctions. *IEEE Transaction on Aerospace and Electronic Systems*, AES-11(4), pp.465–473.

- Clarke, D.W., 1985. Introduction to self-tuning controllers. In C. John Harris & S. A. Billings, eds. *Self-tuning and Adaptive Control: Theory and Applications*. Stevenage, UK: Peter Peregrinus, Ltd., p. 362.
- Cohen, J. and Cohen, P., 2003. *Applied multiple regression/correlation analysis for the behavioral sciences*, Mahwah, NJ: L. Erlbaum Associates.
- Crosier, R.B., 1986. A New Two-Sided Cumulative Sum Quality Control Scheme. *Technometrics*, 28(3), pp.187–194.
- CSU, 2012. Building Metering Guide. The California State University Office of the Chancellor. Available at: http://www.calstate.edu/cpdc/ae/gsf/documents/CSU_Metering_Guide.pdf
- Curtin, J.M., 2007. *The development and testing of an automated building commissioning analysis tool (ABCAT)*. Master's Thesis, Texas A&M University, College Station, TX
- Deng, S., 1997. *Development and application of a procedure to estimate overall building and ventilation parameters from monitored commercial building energy use*. Master's Thesis, Texas A&M University, College Station, TX. Available at: <http://hdl.handle.net/1969.1/ETD-TAMU-1997-THESIS-D46>
- Deru, M., Field, K., Studer, D., Liu, B., Halverson, M., Winiarski, D., Yazdanian, M., Huang, J., and Crawley, D., 2011. *U.S. Department of Energy Commercial Reference Building Models of the National Building Stock*, Golden, CO: National Renewable Energy Laboratory.
- Dhar, A., 1995. *Development of Fourier Series and Artificial Neural Network Approaches to Model Hourly Energy Use in Commercial Buildings*. Texas A&M University, College Station, TX. Available at: <http://hdl.handle.net/1969.1/154825>
- Dhar, A., Reddy, T.A. and Claridge, D.E., 1999. Generalization of the Fourier Series Approach to Model Hourly Energy Use in Commercial Buildings. *Journal of Solar Energy Engineering*, 121(1), pp.54–62.
- Dodier, R.H. and Kreider, J.F., 1999. Detecting whole building energy problems. *ASHRAE Transactions*, 105, p.579.
- DOE, 2006. *Guidance for Electric Metering in Federal Buildings* (Report No. DOE/EE-0312). Energy Efficiency and Renewable Energy, U.S. Department of Energy. Available at: https://www1.eere.energy.gov/femp/pdfs/adv_metering.pdf.
- DOE, 1997. *International Performance Measurement and Verification Protocol (IPMVP)*, Washington, D.C.: U.S. Department of Energy.

- Dong, B., Cao, C. and Lee, S.E., 2005. Applying support vector machines to predict building energy consumption in tropical region. *Energy and Buildings*, 37(5), pp.545–553.
- Effinger, M., Friedman, H. and Moser, D., 2010. Building Performance Tracking in Large Commercial Buildings: Tools and Strategies. Subtask 4.2 Research Report: Investigate Energy Performance Tracking Strategies in the Market. September 2010: California Commissioning Collaborative.
- EISA, 2007. *Energy Independence and Security Act of 2007*, Public Law 100-140, as amended, Section 434(b).
- Emmerich, S.T., Persily, A.K. and McDowell, T.P., 2005. Impact of commercial building infiltration on heating and cooling loads in US office buildings. In *26th AIVC Conference “Ventilation in relation to the energy performance of buildings.”* Brussels, Belgium.
- EO 13514, 2009. *Executive Order (EO) 13514: Federal Leadership in Environmental, Energy, and Economic Performance*, Signed October 8, 2009.
- EPAct, 2005. *Energy Policy Act of 2005*, Public Law 109-58, as amended, Section 103, Energy Use Measurement and Accountability, Section 543 (42 USC 8253), (e) Metering of Energy Use.
- Erpelding, B., 2008. Monitoring chiller plant performance. *ASHRAE Journal*, (April 2008), pp.48–52.
- EVO, 2012. *International Performance Measurement & Verification Protocol: Concept and Options for Determining Energy and Water Saving Volumes I*, The Efficiency Valuation Organization (EVO). Available at: <http://www.evo-world.org/>
- Fels, M.F., 1986. PRISM: An Introduction. *Energy and Buildings*, 9(1–2), pp.5–18.
- Flouquet, F., 1992. Local Weather Correlations and Bias in Building Parameter Estimates from Energy-Signature Models. *Energy and Buildings*, 19(2), pp.113–123.
- Fox, J., 2008. *Applied regression analysis and generalized linear models*, Los Angeles, CA: SAGE Publications.
- Fox, K. and Morton, B., 2013. *World Green Building Trends: Business Benefits Driving New and Retrofit Market Opportunities in Over 60 Countries*, McGraw-Hill Construction. Available at: <http://analyticsstore.construction.com/index.php/world-green-building-trends-smartmarket-report-2013.html>

- Friedman, H., Crowe, E., Sibley, E., and Effinger, M., 2011. *The Building Performance Tracking Handbook*, Available at: <http://www.cacx.org/PIER/documents/bpt-handbook.pdf>.
- Friedman, H. and Piette, M.-A., 2001. Comparative guide to emerging diagnostic tools for large commercial HVAC systems (Report No. LBNL-2899E). Lawrence Berkeley National Laboratory. Available at: <http://escholarship.org/uc/item/37b942xx.pdf>
- Gertler, J., 1991. Analytical Redundancy Methods in Fault Detection and Isolation. In *IFAC/IAMCS symposium on safe process*. Baden-Baden, Germany.
- Granderson, J., Piette, M., Ghatikar, G., and Price, P., 2009. *Building Energy Information Systems: State of Technology and User Case Studies* (Report No. LBNL-2899E). Lawrence Berkeley National Laboratory.
- Gustafsson, F., 2001. *Adaptive Filtering and Change Detection*, West Sussex, England: John Wiley & Sons.
- Haan, C.T., 2002. *Statistical Methods in Hydrology*. Hoboken, NJ: Wiley-Blackwell.
- Haberl, J.S. et al., 1988. An expert system for building energy consumption analysis: applications at a university campus. *ASHRAE Transactions*, 94(1), pp.1037–1062.
- Haberl, J.S., Norford, L., Spadaro, J., 1989. Diagnosing Building Operational Problems: Intelligent Systems for Diagnosing Operational Problems in HVAC Systems, *ASHRAE Journal*, vol.31, no.6, pp.20–30.
- Haberl, J.S. and Claridge, D.E., 1987. An Expert System for Building Energy Consumption Analysis: Prototype Results. *ASHRAE Transactions*, 93(1), pp.979–998.
- Haberl, J.S. and Thamilsaran, S., 1996. Great Energy Predictor Shootout II: Measuring retrofit savings - overview and discussion of results. *ASHRAE Transactions*, 102(2), pp.419–435.
- Hamby, D.M., 1994. A review of techniques for parameter sensitivity analysis of environmental models. *Environmental monitoring and assessment*, 32(2), pp.135–154. Available at: <http://dx.doi.org/10.1007/BF00547132>.
- Hammarsten, S., 1987. A Critical Appraisal of Energy-Signature Models. *Applied Energy*, 26(2), pp.97–110.
- Hammarsten, S., 1984. *Estimation of energy balances for houses* (Bulletin M84:18). National Swedish Institute for Building Research.

- Harris, P., 1989. *Energy Monitoring and Target Setting Using CUSUM*, Cambridge, UK: Cheriton Technology Management, Ltd.
- Haykin, S.S., 2001. *Adaptive filter theory* 4th ed. Englewood Cliffs, NJ: Prentice Hall.
- Henze, G., Kalz, D., Felsmann, C., and Knebe, G., 2004. Impact of forecasting accuracy on predictive optimal control of active and passive building thermal storage inventory. *HVAC&R Research*, 10(2), pp.153–179.
- Hilliard, A. and Jamieson, G., 2013. Recursive Estimates as an Extension to CUSUM-based Energy Monitoring & Targeting. In *The 10th biennial Summer Study on Energy Efficiency in Industry*. Niagara Falls, NY: American Council for an Energy-Efficient Economy.
- Hinkley, D. V, 1971. Inference about the Change-Point from Cumulative Sum Tests. *Biometrika*, 58(3), pp.509–523.
- Hyvarinen, J. ed., 1996. *IEA Annex 25 Final Report Vol 1: Real Time Simulation of HVAC Systems for Building Optimizaation, Fault Detection and Diagnosis SourceBook*, International Energy Agency.
- Hyvarinen, J. ed., 1997. *IEA Annex 25 Final Report Vol 2: Technical Papers of IEA Annex 25*. International Energy Agency.
- ICC, 2009. *2009 International Energy Conservation Code*, Country Club Hills, IL: International Code Council, Inc.
- Isermann, R., 1997. Supervision, fault-detection and fault-diagnosis methods — An introduction. *Control Engineering Practice*, 5(5), pp.639–652.
- Jagpal, R. ed., 2006. *Technical Synthesis Report: Computer Aided Evaluation of HVAC System Performance*. Hertfordshire, UK: Faber Maunsell Ltd.
- JCGM, 2010. *Evaluation of measurement data – Guide to the expression of uncertainty in measurement* (JCGM 100:2008). The Joint Committee for Guides in Metrology. Available at: www.bipm.org.
- Ji, J., Baltazar, J.C. and Claridge, D.E., 2008. Study of the Outside Air Enthalpy Effects in the Screening of Metered Building Energy Data. *The 16th symposium on improving building systems in hot and humid climates*, December 2008. Plano, TX: Energy Systems Laboratory, Texas A&M University.
- Jones, B. and Nachtsheim, C., 2011. A Class of Three-Level Designs for Definitive Screening in the Presence of Second-Order Effects. *Journal of Quality Technology*, 43, pp.1–15.

- Karatasou, S., Santamouris, M. and Geros, V., 2006. Modeling and predicting building's energy use with artificial neural networks: Methods and results. *Energy and Buildings*, 38(8), pp.949–958.
- Katipamula, S. and Brambley, M.R., 2005. Methods for Fault Detection, Diagnostics, and Prognostics for Building Systems-- A Review, Part I. *HVAC&R Research*, 11(1), pp.169–187.
- Katipamula, S., Reddy, T.A. and Claridge, D.E., 1994. Development and Application of Regression Models to Predict Cooling Energy Consumption in Large Commercial Buildings. In *the 1994 ASME/JSME/JSES International Solar Energy Conference*. San Francisco, California: The American Society of Mechanical Engineers, pp. 307–322.
- Katipamula, S., Reddy, T.A. and Claridge, D.E., 1998. Multivariate Regression Modeling. *Journal of Solar Energy Engineering*, 120(3), pp.177–184.
- Kim, K.H., 2014. *Development of an Improved Methodology for Analyzing Existing Single-Family Residential Energy Use*. Texas A&M University. Available at: <http://hdl.handle.net/1969.1/153252>.
- Kim, K.H. and Haberl, J.S., 2015. Development of methodology for calibrated simulation in single-family residential buildings using three-parameter change-point regression model. *Energy and Buildings*, 99, pp.140–152. Available at: <http://www.sciencedirect.com/science/article/pii/S037877881500331X> [Accessed January 12, 2016].
- Kim, S., 2012. *ppcor: Partial and Semi-partial (Part) correlation*, Available at: <http://cran.r-project.org/package=ppcor>.
- Kim, W. and Katipamula, S., 2017. A review of fault detection and diagnostics methods for building systems. *Science and Technology for the Built Environment*, DOI: 10.1080/23744731.2017.1318008
- Kissock, J.K., 1993. *A Methodology to Measure Retrofit Energy Savings in Commercial Buildings* (Dissertation No. 9410821). College Station, TX: Texas A&M University.
- Kissock, J.K., Haberl, J.S. and Claridge, D.E., 2004. *ASHRAE 1050-RP: Development of a Toolkit for Calculating Linear, Change-Point Linear and Multiple-Linear Inverse Building Energy Analysis Models*, Atlanta, GA.: American Society of Heating, Refrigerating, and Air-Conditioning Engineers.
- Kissock, J.K., Reddy, T.A. and Claridge, D.E., 1998. Ambient temperature regression analysis for estimating retrofit savings in commercial buildings. *ASME Journal of*

Solar Energy Engineering, 120, pp.168–176.

Kramer, H., Russell, J. and Crowe, E., 2013. *Inventory of Commercial Energy Management and Information Systems (EMIS) for M&V Applications: Final Report* (Report No. E13-264). Portland Energy Conservation, Inc.; North Energy Efficiency Alliance.

Krarti, M., 2012. *Weatherization and Energy Efficiency Improvement for Existing Homes: An Engineering Approach*, Boca Raton, FL: CRC Press.

Kreider, J.F. and Haberl, J.S., 1994. Predicting hourly building energy use: the great energy predictor shootout - overview and discussion of results. *ASHRAE Transactions*, 100(2), pp. 1104–1118.

LBNL, 2011. eQUEST. Available at:
<http://www.doe2.com/download/equest/eQUESTv3-Overview.pdf>.

Leung, S.-H. and So, C.F., 2005. Gradient-based variable forgetting factor RLS algorithm in time-varying environments. *Signal Processing, IEEE Transactions on Signal Processing*, 53(8), pp.3141–3150.

Liddament, M., 1999. *Technical Synthesis Report: Real Time Simulation of HVAC Systems for Building Optimisation, Fault Detection and Diagnostics*, Coventry, UK: ESSU.

Liu, F., Jiang, H., Lee, Y., Snowdon, J., and Bobker, M., 2011. *Statistical Modeling for Anomaly Detection, Forecasting and Root Cause Analysis of Energy Consumption for a Portfolio of Buildings*, IBM Research Report (RC25165W)

Ljung, L., 1999. *System Identification* 2nd ed., Upper Saddle River, NJ: Prentice Hall PTR.

Lowe, R., Wingfield, J., Bell, M., and Bell, J., 2007. Evidence for heat losses via party wall cavities in masonry construction. *Building Services Engineering Research and Technology*, 28(2), pp.161–181. Available at: <http://discovery.ucl.ac.uk/3377/>.

Lucas, J.M. and Crosier, R.B., 2000. Fast initial response for CUSUM quality-control schemes: Give your CUSUM a head start. *Technometrics*, 42(1), p.102.

MacGregor, J.F. and Harris, T.J., 1993. The exponentially weighted moving variance. *Journal of Quality Technology*, 25(2), pp.106–118.

Maile, T., Bazjanac, V. and Fischer, M., 2012. A method to compare simulated and measured data to assess building energy performance. *Building and Environment*, 56, pp.241–251.

- Masuda, H., Baltazar, J.C., Ji, J., and Claridge, D.E., 2008. Development of Data Quality Control Limits for Data Screening through the “Energy Balance” Method. In *Proceedings of the Sixteenth Symposium on Improving Building Systems in Hot and Humid Climates*. Plano, TX: Energy Systems Laboratory, Texas A&M University.
- Masuda, H., Ji, J., Baltazar, J.C., and Claridge, D.E., 2009. Use of First Law Energy Balance as a Screening Tool for Building Energy Use Data: Experiences on the Inclusion of Outside Air Enthalpy Variable. In *Proceedings of the Ninth International Conference for Enhanced Building Operations*. Austin, TX: Energy Systems Laboratory, Texas A&M University.
- Masuda, H. and Claridge, D.E., 2012. Estimation of Building Parameters Using Simplified Energy balance Model and Metered Whole Building Energy Use. In *The Twelfth International Conference for Enhanced Building Operations*. Manchester, UK: Energy Systems Laboratory, Texas A&M University.
- Mathews, J. and Douglas, S., 2001. *Adaptive Filters*, Prentice Hall PTR. Available at: <http://books.google.com/books?id=6iEgAAAACAAJ>.
- Mathieu, J., Price, P., Kiliccote, S., and Piette, M.-A., 2011. Quantifying changes in building electricity use, with application to demand response. *IEEE Transaction on Smart Grid*, 2(3), pp.507-518.
- Mills, E., 2011. Building commissioning: a golden opportunity for reducing energy costs and greenhouse gas emissions in the United States. *Energy Efficiency*, 4(2), pp.145-173.
- Montgomery, D.C., 2009. *Introduction to statistical quality control*, : John Wiley & Sons, c2009.
- Montgomery, D.C., Peck, E.A. and Vining, G.G., 2012. *Introduction to linear regression analysis*, Hoboken, NJ : Wiley.
- Morris, A.S. and Langari, R., 2012. *Measurement and Instrumentation: Theory and Application*. Boston, MA: Butterworth-Heinemann.
- Newsham, G.R., Mancini, S. and Birt, B.J., 2009. Do LEED-certified buildings save energy? Yes, but.... *Energy and Buildings*, 41(8), pp.897–905.
- NOAA, 2012. Quality Controlled Local Climatological Data. *QCLCD*. National Oceanic and Atmospheric Administration. Available at: <https://www.ncdc.noaa.gov/> [Accessed July 1, 2012].
- NOAA, 2013. Quality Controlled Local Climatological Data (2.5.4). *QCLCD (2.5.4)*. National Oceanic and Atmospheric Administration. Available at:

<https://www.ncdc.noaa.gov/> [Accessed January 5, 2013].

- O’Neal, D.L., Bryant, J.A., Turner, W.D., and Glass, M.G., 1990. Metering and Calibration in LoanSTAR Buildings. In *Symposium on Improving Building Systems in Hot and Humid Climates*. Available at: <http://hdl.handle.net/1969.1/6595>.
- O’Neill, Z., Pang, X., Shashanka, M., Haves, P., & Bailey, T. (2014). Model-based real-time whole building energy performance monitoring and diagnostics. *Journal of Building Performance Simulation*, 7(2), 83–99.
- Oates, D. and Sullivan, K.T., 2012. Postoccupancy Energy Consumption Survey of Arizona’s LEED New Construction Population. *Journal of Construction Engineering and Management*, 138(6), pp.742–750.
- Page, E.S., 1954. Continuous Inspection Schemes. *Biometrika*, 41(1/2), pp.100–115.
- Paleologu, C., 2008. A robust variable forgetting factor recursive least-squares algorithm for system identification. *IEEE Signal Processing Letters*, 15(3), pp.597–600.
- Palmiter, L.S., Hamilton, L.B. and Holtz, M.J., 1979. *Low cost performance evaluation of passive solar buildings* (Report No. SERI/RR-63-223). Solar Energy Research Institute, Golden, CO.
- Patton, R.J., Frank, P.M. and Clark, R.N., eds., 2000. *Issues of Fault Diagnosis for Dynamic Systems*. London, UK: Springer.
- Rabl, A., 1988. Parameter Estimation in Buildings: Methods for Dynamic Analysis of Measured Energy Use. *Journal of Solar Energy Engineering-Transactions of the ASME*, 110, pp.52–66.
- Rabl, A. and Rialher, A., 1992. Energy Signature Models for Commercial Buildings - Test with Measured Data and Interpretation. *Energy and Buildings*, 19(2), pp.143–154.
- R Core Team, 2013. *R: A Language and Environment for Statistical Computing*, Vienna, Austria: R Foundation for Statistical Computing. Available at: <http://www.r-project.org/>.
- Ramanathan, R., Engle, R., Granger, C., Vahid-Argahi, F., and Brace, C., 1997. Short-run forecasts of electricity loads and peaks. *International Journal of Forecasting*, 13(2), pp.161–174.
- Reddy, T.A., Kissock, J.K., Katipamula, S., and Claridge, D.E., 1994. An energy delivery efficiency index to evaluate simultaneous heating and cooling effects in large commercial buildings. *Journal of Solar Energy Engineering*, 116(May),

pp.79–87.

- Reddy, T.A., Liu, M., Katipamula, S., and Claridge, D.E., 1998. Extending the Concept of Energy Delivery Efficiency (EDE) of HVAC Systems. *ASHRAE Transactions*, 104, p.13.
- Reddy, T.A., Katipamula, S., Kisoock, J.K., and Claridge, D.E., 1995. The Functional Basis of Steady-State Thermal Energy Use in Air-Side HVAC Equipment. *Journal of Solar Energy Engineering*, 117(1), pp.31–39.
- Reddy, T.A., Deng, S. and Claridge, D.E., 1999. Development of an Inverse Method to Estimate Overall Building and Ventilation Parameters of Large Commercial Buildings. *Journal of Solar Energy Engineering*, 121, pp.40–46.
- Reddy, T.A., Kisoock, J.K. and Ruch, D.K., 1998. Uncertainty in Baseline Regression Modeling and in Determination of Retrofit Savings. *Journal of Solar Energy Engineering*, 120, pp.185–192.
- Reichmuth, H. and Egnor, T., 2013. X-View™: The Reichmuth Framework II. In *Technologies for Sustainability (SusTech), 2013 1st IEEE Conference on*. pp. 258–261.
- Reichmuth, H. and Turner, C., 2010. A Tool for Efficient First Views of Commercial Building Energy Performance. , pp.325–338.
- Richalet, V., Neirac, F.P., Tellez, F., Marco, J., and Bloem, J.J., 2001. HELP (house energy labeling procedure): methodology and present results. *Energy and Buildings*, 33(3), pp.229–233.
- Roberts, S.W., 1959. Control Chart Tests Based on Geometric Moving Averages. *Technometrics*, 1(3), pp.239–250.
- Robinson, J., 1992. Flowloop at the Energy Systems Laboratory. , pp.1–29.
- Ruch, D. and Claridge, D.E., 1992. A Four-Parameter Change-Point Model for Predicting Energy Consumption in Commercial Buildings. *Journal of Solar Energy Engineering*, 114(May), pp.77–83.
- Ruch, D.K., Kisoock, J. and Reddy, T., 1999. Prediction Uncertainty of Linear Building Energy Use Models with Autocorrelated Residuals. *Journal of Solar Energy Engineering*, 121(February), pp.63–68.
- Ruch, D.K., Kisoock, J.K. and Reddy, T.A., 1999. Prediction Uncertainty of Linear Building Energy Use Models with Autocorrelated Residuals. *Journal of solar energy engineering*, 121(1), pp.63–68.

- Ryan, T.P., 2011. *Statistical Methods for Quality Improvement* 3rd ed., San Francisco, CA: John Wiley & Sons, Inc.
- Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Carboni, J., Gatelli, D., Saisana, M., and Tarantola, S., 2008. *Global Sensitivity Analysis. The Primer*, England, UK: John Wiley & Sons, Ltd.
- Saltelli, A., 2002. Making best use of model evaluations to compute sensitivity indices. *Computer Physics Communications*, 145(2), pp.280–297. Available at: <http://www.sciencedirect.com/science/article/pii/S0010465502002801> [Accessed December 2, 2015].
- SAS, 2014. *JMP 11 Profilers*, Cary, NC: SAS Institute Inc.
- Schrock, D.W. and Claridge, D.E., 1989. Predicting Energy Usage in a Supermarket. In *Proceedings of the Sixth Symposium on Improving Building Systems in Hot and Humid Climates*. Dallas, Texas, p. F-19-F-27.
- Schwarz, G., 1978. Estimating the Dimension of a Model. *Annals of Statistics*, 6(2), pp.461–464.
- Shao, X., 2006. *First Law Energy Balance as a Data Screening Tool*. College Station, TX: Texas A&M University.
- Shao, X. and Claridge, D.E., 2006. Use of first law energy balance as a screening tool for building energy data, part I - Methodology. *ASHRAE Transactions*, 112(2), pp.717–731.
- Sjögren, J.-U., Andersson, S., and Olofsson, T., 2007. An Approach to Evaluate the Energy Performance of Buildings Based on Incomplete Monthly Data. *Energy and Buildings*, 39(8), pp.945–953.
- Sjögren, J.-U., Andersson, S. and Olofsson, T., 2009. Sensitivity of the total heat loss coefficient determined by the energy signature approach to different time periods and gained energy. *Energy and Buildings*, 41(7), pp.801–808.
- Solupe, M. and Krarti, M., 2014. Assessment of infiltration heat recovery and its impact on energy consumption for residential buildings. *Energy Conversion and Management*, 78, pp.316–323.
- Somogyi, Z., 1998. *In situ evaluation of the thermal characteristics of building components and buildings including the comparison with predicted performances*. Universite catholique de Louvain.
- Sonderegger, R.C., 1998. A baseline model for utility bill analysis using both weather

- and non-weather-related variables. *ASHRAE Transactions*, 104, p.859.
- Sonderegger, R.C., 1978. Diagnostic Tests Determining the Thermal Response of a House. *ASHRAE Transactions*, 1, pp.691–702.
- Sonderegger, R.C. and Modera, M.P., 1979. Electric Co-heating : A Method for Evaluating Seasonal Heating Efficiencies and Heat Loss Rates in Dwellings. In *The Second CIB Symposium on Energy Conservation in the Built Environment*. Copenhagen, Denmark, May 28-June 1: Lawrence Berkeley National Laboratory.
- Stuart, E., Larsen, P.H., Goldman, C.A., and Gilligan, D., 2013. Current Size and Remaining Market Potential of the U.S. Energy Service Company Industry. Lawrence Berkeley National Laboratory.
- Stuart, G., Fleming, P., Ferreira, V., and Harris, P., 2007. Rapid analysis of time series data to identify changes in electricity consumption patterns in UK secondary schools. *Building and Environment*, 42(4), pp.1568–1580.
- Subbarao, K., 1988. Short-Term Energy Monitoring (STEM): Application of the PSTAR Method to a Residence in Fredericksburg, Virginia (Report No. SERI/TR-254-3356). Solar Energy Research Institute, Golden, CO.
- Team, R.C., 2011. *R: A Language and Environment for Statistical Computing*, Vienna, Austria: R Foundation for Statistical Computing. Available at: <http://www.r-project.org/>.
- UIUC and LBNL, 2007. EnergyPlus Engineering Reference: the reference to EnergyPlus calculations.
- Ulickey, J., Fackler, T., Koeppel, E., and Soper, J., 2010. *Building Performance Tracking in Large Commercial Buildings : Tools and Strategies: Characterization of Fault Detection and Diagnostic (FDD) and Advanced Energy Information Systems (EIS) Tools*, California Commissioning Collaborative. Available at: https://www.cacx.org/PIER/documents/Subtask_4-3_Report.pdf
- USGBC, 2014. *Leadership in Energy & Environmental Design*, The U.S. Green Building Council. Available at: <http://www.usgbc.org/leed>.
- Wald, A., 1947. *Sequential analysis*, NY: Wiley & sons, inc.
- Watt, J. and Haberl, J., 1994. Flow Measurement with Tangential Paddlewheel Flow Meters: Analysis of Experimental Results and In-situ Diagnostics. In *The Sixteenth National Industrial Energy Technology Conference*. Houston, TX: Energy Systems Laboratory, Texas A&M University.

- White, J. and Reichmuth, R., 1996. Simplified method for predicting building energy consumption using average monthly temperatures. *Energy Conversion Engineering Conference, 1996. IECEC 96., Proceedings of the 31st Intersociety*. Washington, D.C., pp. 1834–1839.
- de Wilde, P., Martinez-Ortiz, C., Pearson, D., Beynon, I., Beck, M., and Barlow, N., 2013. Building simulation approaches for the training of automated data analysis tools in building energy management. *Advanced Engineering Informatics*, 27(4), pp.457–465.
- Wingfield, J., Johnston, D., Miles-Shenton, D., and Bell, M., 2010. *Whole House Heat Loss Test Method (Coheating)*. Centre for Built Environment. Available at: [http://www.leedsmet.ac.uk/as/cebe/projects/iea_annex58/whole_house_heat_loss_test_method\(coheating\).pdf](http://www.leedsmet.ac.uk/as/cebe/projects/iea_annex58/whole_house_heat_loss_test_method(coheating).pdf).
- Wortham, A.W., Heinrich, G.F., and Taylor, D., 1974. Adaptive exponentially smoothed control charts. *International Journal of Production Research*, 12(6), pp.683–690.
- Yalcintas, M. and Akkurt, S., 2005. Artificial neural networks applications in building energy predictions and a case study for tropical climates. *International Journal of Energy Research*, 29(10), pp.891–901.
- Yang, J., Rivard, H. and Zmeureanu, R., 2005. On-line building energy prediction using adaptive artificial neural networks. *Energy and Buildings*, 37(12), pp.1250–1259.
- Young, P.C., 2011. *Recursive Estimation and Time-Series Analysis: An Introduction for the Student and Practitioner* 2nd ed., UK: Springer.
- Yu, B. and VanPaassen, D.H.C., 2003. Fuzzy neural networks model for building energy diagnosis. In *Eith International IBPSA Conference*. Eindhoven, Netherlands, pp. 1459–1466.
- Zeileis, A., Kleiber, C., Krämer, W., and Hornik, K., 2003. Testing and dating of structural changes in practice. *Computational Statistics & Data Analysis*, 44(1–2), pp.109–123.
- Zhang, J., Xie, Y., Athalye, R., Goel, S., Hart, R., Mendon, V., Rosenberg, M., and Liu, B., 2013. *Energy and Energy Cost Savings Analysis of the IECC for Commercial Buildings*, Richland, Washington. August, 2013. Available at: <http://www.energycodes.gov/sites/default/files/documents/PNNL-22760.pdf>.
- Zoebelein, M., 2014. Flow Research: Magnetic Now the Largest Flowmeter Market. Available at: <http://www.flowcontrolnetwork.com/flow-research-magnetic-now-the-largest-flowmeter-market/> [Accessed January 1, 2016].

APPENDIX A

ENERGY USE AND E_{BL} DATA FOR REGRESSION ANALYSIS IN CHAPTER 2

Table A-1 shows the building numbers and corresponding data filenames. Each file contains the data listed in Table A-2.

Table A-1. Building number and filename	
Building number	Filename
1	0290_Wells
2	0291_Rudder
3	0292_Eppright
4	0293_Appelt
5	0361_Bright Football Complex
6	0369_Read and GRWA
7	0383_Koldus
8	0384_Sanders
9	0386_Jack E Brown
10	0398_Langford A
11	0400_Spence
12	0403_Fountain
13	0412_Moses
14	0415_DavisGary
15	0419_Legett
16	0420_Milner
17	0426_FHK
18	0430_Schumacher
19	0435_Harrington EC
20	0439_Cain

Table A-1 Continued.

Building number	Filename
21	0446_Rudder Theatre Complex
22	0449_BSBW
23	0462_Academic
24	0465_Butler
25	0470_Glasscock
26	0471_Pavilion
27	0473_Williams Admin
28	0480_Bolton
29	0482_Fermier
30	0490_Halbouty
31	0507_Vet Med Science
32	0518_Zachry
33	0520_Beutel
34	0521_Heldenfels
35	0524_Blocker
36	0548_Clements
37	0549_Haas
38	0652_Neeley
39	0972_Lab Animal Care
40	1085_Vet SAH
41	1184_Bio Waste
42	1194_Vet LA
43	1497_UBO
44	1502_Heep Center
45	1504_Raynolds

Table A-1 Continued.	
Building number	Filename
46	1506_Horticulture
47	1508_Price Hobgood
48	1511_West Campus Library
49	1525_Nuclear Magnetic
50	1530_Interdisciplinary
51	1560_Student Rec Center
52	1606_Bush Library
53	1607_Allen
54	1609_TTI Headquarters
55	1810_State Chemist
56	1904_TIPS A

Table A-2. Data description		
Column Title	Description	Unit
Date	Date	
DryBulbFahrenheit	Daily average outside air temperature	°F
DryBulbCelsius	Daily average outside air temperature	°C
WetBulbFahrenheit	Daily average wetbulb temperature	°F
WetBulbCelsius	Daily average wetbulb temperature	°C
DewPointFahrenheit	Daily average dewpoint temperature	°F
DewPointCelsius	Daily average dewpoint temperature	°C
RelativeHumidity	Daily average relative humidity	% r.h.
StationPressure	Daily average station pressure	in.Hg
Humidity Ratio	Daily average humidity ratio	lb/lb _{da} or kg/kg _{da}

Table A-2 Continued.

Column Title	Description	Unit
EnthalpyIP	Daily average enthalpy	kJ/kg _{da}
EnthalpySI	Daily average enthalpy	Btu/lb _{da}
HRLoad	W_{oa}^+	lb/lb _{da} or kg/kg _{da}
Af_sqft	Floor area	ft ²
Af_sqm	Floor area	m ²
ELE_kWh	Daily electricity use	kWh
CHW_MMBtu	Daily chilled water use	MMBtu
HHW_MMBtu	Daily heating hot water use	MMBtu
ELE_MMBtu	Daily electricity use	MMBtu
EBL_MMBtu	E_{BL}	MMBtu
EBL_Btu_sqft	E_{BL}	Btu/ft ²
ELE_GJ	Daily electricity use	GJ
CHW_GJ	Daily chilled water use	GJ
HHW_GJ	Daily heating hot water use	GJ
EBL_GJ	E_{BL}	GJ
EBL_kJ_sqm	E_{BL}	kJ/m ²
Solar_Wm2	Daily average solar insolation	W/m ²

APPENDIX B

ENERGYPLUS INPUT FILES FOR CHAPTER 3

B.1 As-is Case (AsIs.idf)

B.2 Ideal Case (Ideal.idf)

APPENDIX C

ENERGY USE, E_{BL} AND Q_B DATA FOR REGRESSION ANALYSIS IN CHAPTER 3

C.1 As is

C.2 Ideal w/ solar

C.3 Ideal w/o solar

C.4 Haas Hall

C.5 McFadden Hall

C.6 Hobby Hall

APPENDIX D

FACTOR LEVELS AND E_{BL} PARAMETER ESTIMATES FOR SIMULATION RUNS

IN CHAPTER 4

Table D-3. The levels of continuous factors in the simulation runs

Order	Bldg. Azimuth (degrees)	Floors	Absorptance	Indoor T (°F)	Occupancy (ft ² /person)	Ventilation (cfm/ft ²)	Wall R Btu (°F·ft ² ·hr/Btu)	Window ratio
1	-45	7	0.8	76	320	0.18	8.13	0.5
2	-45	1	0.4	72	120	0.06	5.65	0.2
3	0	4	0.4	76	120	0.06	5.65	0.5
4	-90	4	0.8	72	320	0.18	8.13	0.2
5	0	7	0.6	72	320	0.06	5.65	0.2
6	-90	1	0.6	76	120	0.18	8.13	0.5
7	0	1	0.8	74	120	0.18	5.65	0.2
8	-90	7	0.4	74	320	0.06	8.13	0.5
9	0	7	0.4	76	220	0.06	8.13	0.2
10	-90	1	0.8	72	220	0.18	5.65	0.5
11	0	7	0.8	72	320	0.12	5.65	0.5
12	-90	1	0.4	76	120	0.12	8.13	0.2
13	0	7	0.8	76	120	0.18	6.89	0.2
14	-90	1	0.4	72	320	0.06	6.89	0.5
15	0	1	0.8	76	320	0.06	8.13	0.35
16	-90	7	0.4	72	120	0.18	5.65	0.35
17	0	1	0.4	76	320	0.18	5.65	0.5
18	-90	7	0.8	72	120	0.06	8.13	0.2
19	0	1	0.4	72	320	0.18	8.13	0.2
20	-90	7	0.8	76	120	0.06	5.65	0.5
21	0	7	0.4	72	120	0.18	8.13	0.5
22	-90	1	0.8	76	320	0.06	5.65	0.2
23	0	1	0.8	72	120	0.06	8.13	0.5
24	-90	7	0.4	76	320	0.18	5.65	0.2
25	-45	4	0.6	74	220	0.12	6.89	0.35
26	-45	4	0.6	74	220	0.12	6.89	0.35

Table D-4. The levels of categorical factors in the simulation runs

Order	Aspect ratio	Glass type	Cold deck reset
1	rectangle	double	on
2	square	single	off
3	rectangle	double	off
4	square	single	on
5	rectangle	double	on
6	square	single	off
7	square	double	on
8	rectangle	single	off
9	square	single	on
10	rectangle	double	off
11	square	single	off
12	rectangle	double	on
13	rectangle	single	off
14	square	double	on
15	square	double	off
16	rectangle	single	on
17	rectangle	single	on
18	square	double	off
19	rectangle	double	off
20	square	single	on
21	square	double	on
22	rectangle	single	off
23	rectangle	single	on
24	square	double	off
25	square	single	off
26	rectangle	double	on

Table D-5. The parameter estimates of the E_{BL} regression model for each simulation run

Order	Standard Error of						Degrees of Freedom	RMSE	Adj.R
	Parameter Estimates			Estimates					
	Intercept $\hat{\beta}_0$	T_{oa} $\hat{\beta}_r$	W_{oa}^+ $\hat{\beta}_w$	Intercept $\hat{\beta}_0$	T_{oa} $\hat{\beta}_r$	W_{oa}^+ $\hat{\beta}_w$			
1	521.9	-7.77	-18437	4.3	0.074	307	362	10.5	0.996
2	341.6	-5.63	-7594	4.2	0.073	302	362	10.3	0.989
3	272.7	-4.60	-6334	2.8	0.048	197	362	6.7	0.993
4	488.9	-7.50	-18874	3.4	0.058	242	362	8.3	0.997
5	241.2	-3.90	-6768	2.2	0.038	157	362	5.3	0.994
6	614.0	-9.29	-18965	4.7	0.081	338	362	11.6	0.996
7	584.5	-9.11	-20510	6.3	0.109	453	362	15.5	0.993
8	359.0	-5.57	-5715	3.0	0.051	213	362	7.3	0.993
9	244.6	-3.80	-6310	1.7	0.030	124	362	4.2	0.996
10	636.6	-10.32	-19193	7.6	0.131	543	362	18.6	0.991
11	447.1	-7.03	-12495	3.3	0.057	235	362	8.0	0.996
12	409.6	-6.38	-13164	3.8	0.066	273	362	9.3	0.994
13	535.9	-8.15	-18572	4.5	0.077	318	362	10.9	0.995
14	301.4	-4.80	-7480	3.7	0.063	264	362	9.0	0.989
15	271.7	-4.27	-6986	3.2	0.055	227	362	7.7	0.990
16	618.2	-9.78	-18813	4.3	0.074	309	362	10.6	0.997
17	698.9	-10.02	-20517	6.6	0.114	474	362	16.2	0.993
18	194.1	-3.65	-6086	2.0	0.035	143	362	4.9	0.995
19	514.7	-7.92	-19185	4.5	0.078	324	362	11.1	0.995
20	326.4	-5.36	-5748	2.7	0.046	192	362	6.5	0.994
21	465.8	-7.58	-18811	3.8	0.065	272	362	9.3	0.996
22	407.6	-6.25	-7270	6.0	0.103	429	362	14.6	0.980
23	353.8	-6.14	-6190	5.1	0.087	363	362	12.4	0.984
24	509.1	-7.43	-18298	4.1	0.070	291	362	10.0	0.996
25	428.7	-6.62	-12560	3.2	0.054	226	362	7.7	0.996
26	393.0	-6.15	-12822	3.2	0.055	229	362	7.8	0.996

Table D-5 Continued.

Order	T_{ref} (°F)	Mean E_{BL} (Btu/day/ft ²)	Maximum E_{BL} (Btu/day/ft ²)	Minimum E_{BL} (Btu/day/ft ²)
1	67.18	-55.1	256.1	-308.9
2	60.65	-59.9	159.1	-204.2
3	59.30	-55.5	122.6	-173.0
4	65.23	-71.0	242.0	-319.1
5	61.90	-41.0	111.6	-149.2
6	66.09	-67.1	309.5	-343.6
7	64.20	-88.6	275.6	-380.3
8	64.42	-33.1	180.1	-167.7
9	64.39	-29.7	122.2	-132.8
10	61.70	-114.4	289.4	-420.0
11	63.56	-63.2	220.6	-261.7
12	64.15	-58.9	193.3	-248.5
13	65.78	-67.0	270.0	-319.0
14	62.79	-43.8	137.7	-175.2
15	63.61	-36.4	126.2	-157.3
16	63.20	-95.6	291.7	-383.6
17	69.76	-35.8	375.3	-330.9
18	53.13	-69.8	72.7	-169.8
19	64.99	-74.7	255.0	-325.0
20	60.88	-51.5	155.6	-180.3
21	61.42	-99.9	208.0	-350.2
22	65.23	-34.6	200.7	-205.2
23	57.59	-78.1	162.0	-239.0
24	68.55	-44.5	265.8	-283.2
25	64.72	-54.1	214.9	-244.7
26	63.91	-58.6	181.9	-245.8

APPENDIX E

EQUEST INPUT FILES FOR PARAMETRIC SIMULATION RUNS IN CHAPTER 4

Table E-6. Simulation run numbers and input files

Simulation run number	Input file name
1	1.inp
2	2.inp
3	3.inp
4	4.inp
5	5.inp
6	6.inp
7	7.inp
8	8.inp
9	9.inp
10	10.inp
11	11.inp
12	12.inp
13	13.inp
14	14.inp
15	15.inp
16	16.inp
17	17.inp
18	18.inp
19	19.inp
20	20.inp
21	21.inp
22	22.inp
23	23.inp
24	24.inp
25	25.inp
26	26.inp

APPENDIX F

RECURSIVE LEAST SQUARES AND CUSUM CHANGE DETECTION PROGRAM

IN CHAPTER 5

This R script generates time series plots for RLS and MLR estimates for comparison. For the RLS estimates, the data for 13 months are used with the first month as a learning period. The RLS and MLR estimates for 12 months are compared in the graphs.

```
# Inputs
# Energy consumption and weather data for a user defined period

# Outputs
# RLS estimates for Intercept, Tslope, and dWslope plus Tref estimate
# Plots for the estimates and original data (3 years, annual, and monthly)

library(grid)

# directories
#wd <- 'C:/R_Work/RLS_Work'
#dd <- 'C:/R_Work/csvOut'
wd <- getwd()
dd <- './csvOut'
weatherfilename <- 'Weather.csv'

source('plotProperties.R')
source('plotScripts.R')
source('RLS_functions.R')

# List of the data files
setwd(dd)
dataFiles <- list.files(pattern = '*.\\*.dat')
setwd(wd)

# Weather data
weather <- read.table(weatherfilename, header = T, sep = ',', na.strings = "-99")
weather$Date <- as.Date(weather$Date, format = "%m/%d/%Y")
dW <- weather$HumidityRatio - dWthreshold
dW[dW < 0] <- 0
```

```

weather <- cbind(weather, dW)

# data periods
ST <- '2013-03-01'
ED <- '2014-03-31'
ST_lm <- '2013-04-01'
ED_lm <- '2014-03-31'
ST_pred <- '2013-04-01'
ED_pred <- '2014-03-31'

# Parameters for RLS
ini.learning <- 10 # used for y ranges in plots
FF_max <- 0.90 # forgetting factor for RLS

# Main process starts from here
for (each.file in dataFiles) { #building loop
  if (regexpr('Weather', each.file) > 0) {
    next # The weather file will be skipped.
  }
  dataTable <- read.table (paste(dd, each.file, sep = '/'), header = T, sep = ',', na.strings = "-99")
  fileName <- unlist (strsplit (each.file, "\\..dat"))[1]
  dataTable$Date <- as.Date(dataTable$Date, format = "%m/%d/%Y", tz="")
  # remove negative consumption values due to cumulative wind-up
  dataTable$ELE_ori[dataTable$ELE_ori < 0] <- NA
  dataTable$CHW_ori[dataTable$CHW_ori < 0] <- NA
  dataTable$HHW_ori[dataTable$HHW_ori < 0] <- NA
  # add EBL to the data table and make dataTable including dataset
  EBL_ori <- calcEBL(dataTable$ELE_ori, dataTable$CHW_ori, dataTable$HHW_ori,
dataTable$Area)
  EBL_mod <- calcEBL(dataTable$ELE_mod, dataTable$CHW_mod, dataTable$HHW_mod,
dataTable$Area)
  EBL_scr <- calcEBL(dataTable$ELE_scr, dataTable$CHW_scr, dataTable$HHW_scr,
dataTable$Area)
  dataTable <- cbind(dataTable, EBL_ori, EBL_mod, EBL_scr)
  #=====
  # Multiple Linear Regression Estimates
  #=====
  cons_lm <- selectTime(dataTable, ST_lm, ED_lm)
  weather_lm <- selectTime(weather, ST_lm, ED_lm)
  cons_lm <- merge(cons_lm, weather_lm, by = 'Date')
  weather_lm_pred <- selectTime(weather, ST_pred, ED_pred)
  cons_pred <- selectTime(dataTable, ST_pred, ED_pred)
  # skip estimation of lm and RLS if all the EBL data are missing
  if ((nrow(cons_lm) == 0) || (nrow(cons_lm[is.na(cons_lm$EBL_ori),]) + 3 >
as.numeric(tail(cons_lm$Date,1) - head(cons_lm$Date,1)))) {
    next
  } else {
    # estimate MLR model
    lm_out <- NULL
    lm_out <- estlm(cons_lm, EBL_ori ~ DryBulbFarenheit + dW, weather_lm_pred,
'DryBulbFarenheit', 'dW')

    #=====

```

```

# Recursive Least Squares Estimates
#=====
consData <- selectTime(dataTable, ST, ED)
weatherData <- selectTime(weather, ST, ED)
consData <- merge(consData, weatherData, by = 'Date')
datalength <- nrow(consData)
InpData <- cbind(consData$DryBulbFahrenheit, consData$dW)#Input signal
OutData <- consData$EBL_ori      #Desired signal

# Initialization
EstOut <- NULL          # Output vector or data frame
h_out <- NULL           # Output of parameter estimates
err <- NA               # filter error
count.NA <- 0          # count NAs from the first day
CL <- NA               # control limit (abs)
alarm <- 0             # abs(error) > control limit, alarm = 1, otherwise 0
errNumber <- 0
varTimer <- 0
# CUSUM
g_1 <- 0
g_2 <- 0
nu <- 0.5
h <- 5
NumDays_1stMonth <- 31
var_v <- 0             # initial value of recursive variance update
# RLS
h_n <- matrix(c(0, 0, 0), nrow = 3, ncol = 1) # coefficient vector
P_n <- 1000000 * diag(nrow(h_n))
k <- 0
FF <- FF_max
EW_var <- FF
var_1 <- NA
var_2 <- 100

# recursive estimations for each building
for (n in 1:datalength){
  x_n <- matrix(c(InpData[n, 1], InpData[n, 2], 1), nrow = 1, ncol = 3)
  d_n <- OutData[[n]]
  varTimer <- varTimer + 1

  if (is.na(d_n[[1]])) {
    count.NA <- count.NA + 1
  }
  if (is.na(d_n) | is.na(x_n[1,1]) | is.na(x_n[1,2])) {
    # If there are any missing data in the input and output
    # the values from the previous step are used.
    h_n <- h_n
    Tref <- NA
    h_out <- c(NA, NA, NA)
    alarm <- 0
    FF <- FF
    err_pri <- NA
    err_pos <- NA
  }
}

```

```

g_1_out <- NA
g_2_out <- NA
} else {
  # k = 1 if CUSUM > h in the previous step
  k = k + 1
  # forget the poor initial estimates but progressively utilize
  # later data with equal weighting
  # Young (4.17)
  FF <- 1 - (1 - FF_max)/(1 - FF_max^(k + 1))
  EW_var <- 1/FF * (EW_var - EW_var^2 / (FF + EW_var))
  # RLS update
  StepEst <- estRLS(x_n, d_n, h_n, P_n, FF)
  h_n <- StepEst$parest
  P_n <- StepEst$invcor
  err_pri <- StepEst$err_pri # priori error = e_n - pred(y)_n-1
  err_pos <- StepEst$err_pos
  # posteriori variance
  var_1 <- StepEst$var_pos
  # Cusum statistics
  if (n > NumDays_1stMonth) {
    norm_err <- err_pri / sqrt(var_2)
    g_1 <- max(g_1 + norm_err - nu, 0)
    g_2 <- max(g_2 - norm_err - nu, 0)
  }
  Tref <- h_n[3, 1] / h_n[1, 1] * -1
  h_out <- c(h_n[1,1], h_n[2,1], h_n[3,1])
  g_1_out <- g_1
  g_2_out <- g_2
  # empirical variance
  var_2 <- var_2 + EW_var * ((err_pri)^2 - var_2)
  # CUSUM detection
  if ((n > NumDays_1stMonth) && ((g_1 > h) || (g_2 > h))) {
    alarm <- 1
    # reset statistics and parameters
    g_1 <- 0
    g_2 <- 0
    k <- 0 # for FF reset
    var_2 <- var_2
    varTimer <- 0
  } else {
    alarm <- 0
  }
}

EstOut <- rbind(EstOut, c(consData$Date[[n]], x_n[1,1], x_n[1,2], d_n[[1]], h_out[1],
h_out[2], h_out[3], err_pri, err_pos, Tref, alarm, FF, EW_var, g_1_out, g_2_out, h, var_1, var_2))
}

EstOut <- data.frame(EstOut)
colnames(EstOut) <- c('Date', 'TDB', 'dW', 'EB', 'Est_T', 'Est_dW', 'Est_Int', 'Pri_Err',
'Post_Error', 'T_ref', 'alarm', 'FF', 'EW_var', 'CUSUM_g_1', 'CUSUM_g_2', 'Threshold', 'ExactVar',
'EmpiricalVar')
EstOut$Date <- as.Date(EstOut$Date, origin = '1970-01-01')

```

```

outfileName <- sprintf ("%s.csv", fileName)
write.table(EstOut, outfileName, quote = F, col.names = T, row.names = F, sep = ',')

outfileName <- sprintf ("%s_lmsummary.txt", fileName)
sink(outfileName)
      print(lm_out$lm_summary)
sink()

# creating grid plots
#graphFileName <- sprintf ("%s_res.eps", fileName)
graphFileName <- sprintf ("%s_1res.pdf", fileName)
pdf(graphFileName, width = 8.5 , height = 11)
#postscript(graphFileName, width = 8.5 , height = 11)      # Lettersize PDF
#par(oma = c(20, 17, 10, 17))      # outer margin for journal papers
par(oma = c(3, 4, 4, 4))
par(mfcol = c(5, 1))      # plotting direction = row
par(ps = 10)      # basic font size

plotGraph1()
plotGraph_MLR_res()
plotGraph_RLS_res()
plotGraph_TestStat()

# main title of the page - building number and name
mtext(side = 3, line = 1, outer = T, text = fileName, cex = 1)

dev.off()

#graphFileName <- sprintf ("%s_par.eps", fileName)
graphFileName <- sprintf ("%s_2par.pdf", fileName)
pdf(graphFileName, width = 8.5 , height = 11)
#postscript(graphFileName, width = 8.5 , height = 11)
#par(oma = c(25, 17, 10, 17))      # outer margin for journal papers
par(oma = c(3, 4, 4, 4))
par(mfcol = c(5, 1))      # plotting direction = row
par(ps = 10)      # (for a journal paper, ps=8 may be good.)

plotGraph_Int()
plotGraph_Tslope()
plotGraph_W()
plotGraph_Tref()
plotGraph_Var()

# main title of the page - building number and name
mtext(side = 3, line = 1, outer = T, text = fileName, cex = 1)

dev.off()

}      # closing bracket for skipping buildings with missing data for the whole modeling period
}

```

This is a function to estimate and output the regression parameters

```

estlm <- function(.basedata, .lmstructure, .newdata, .var1name, .var2name)
{
  EBL_n <- nrow(.basedata) - sum(is.na(.basedata$EBL_ori))
  if (EBL_n > 20)
  {
    lm_model <- lm (.lmstructure, data = .basedata)
    lm_summary <- summary (lm_model)
    maxEBL <- .basedata$EBL_ori[which.max(.basedata$EBL_ori)]
    minEBL <- .basedata$EBL_ori[which.min(.basedata$EBL_ori)]
    Est.p1 <- lm_summary$coefficients ["(Intercept)", "Estimate"]
    Est.p2 <- lm_summary$coefficients [.var1name, "Estimate"]
    Est.p3 <- lm_summary$coefficients [.var2name, "Estimate"]

    # note: errors and p-values may not be returned by lm if the newdata is not NA.
    StdErr.p1 <- lm_summary$coefficients ["(Intercept)", "Std. Error"]
    StdErr.p2 <- lm_summary$coefficients [.var1name, "Std. Error"]
    StdErr.p3 <- lm_summary$coefficients [.var2name, "Std. Error"]
    pVal.p1 <- lm_summary$coefficients ["(Intercept)", "Pr(>|t|)"]
    pVal.p2 <- lm_summary$coefficients [.var1name, "Pr(>|t|)"]
    pVal.p3 <- lm_summary$coefficients [.var2name, "Pr(>|t|)"]
    df <- lm_model$df.residual
    RMSE <- lm_summary$sigma
    AdjR <- lm_summary$adj.r.squared
    CV <- RMSE/((maxEBL - minEBL)/2)
    T_EB0 <- -1 * Est.p1 / Est.p2
    T_EB0_SE <- (Est.p1/Est.p2)*sqrt((StdErr.p1/Est.p1)^2 + (StdErr.p2/Est.p2)^2)

    pred <- predict (lm_model, .newdata, level = CONF.LEVEL)
    predInt <- predict (lm_model, .newdata, level = CONF.LEVEL, interval="prediction")
    confInt <- predict (lm_model, .newdata, level = CONF.LEVEL, interval="confidence")
    colnames(predInt) <- c("pred.fit", "pred.lwr", "pred.upr")
    colnames(confInt) <- c("conf.fit", "conf.lwr", "conf.upr")
    outdata <- cbind(.newdata, predInt, confInt)

    return (list(Est_int = Est.p1, Est_T = Est.p2, Est_W = Est.p3,
                 SE_int = StdErr.p1, SE_T = StdErr.p2, SE_W = StdErr.p3,
                 p_int = pVal.p1, p_T = pVal.p2, p_W = pVal.p3,
                 df = df, RMSE = RMSE, AdjR = AdjR, CV = CV,
                 outdata = outdata, lm_summary = lm_summary, T_EB0 = T_EB0,
                 T_EB0_SE = T_EB0_SE))
  } else {
    return (as.list(rep(NA, 17)))
  }
}

# function to calculate EBL from energy use data
calcEBL <- function(.ELEdata, .CHWdata, .HHWdata, .Area)
{
  return ((.ELEdata * 3412.14 + .HHWdata * 1000000 - .CHWdata * 1000000) / .Area)
}

# function to extract specified period from a dataframe
selectTime <- function(.data, .startdate, .enddate)
{

```



```

        return(.data[.data$Date >= as.Date(.startdate) & .data$Date <= as.Date(.enddate),])
    }

# function to update RLS equations at each recursion
estRLS <- function(x, d, h, P, FF)
{
    # priori error (one-step ahead prediction error)
    e_pri <- drop(d - x %*% h)
    # gain vector
    k <- P %*% t(x) / drop(FF + x %*% P %*% t(x))
    # new parameter estimates
    h <- h + k * e_pri
    # new P
    P <- (P - k %*% x %*% P) / FF
    # posteriori error
    e_pos <- drop(d - x %*% h)
    # error variance
    v_pos <- 1 + x %*% P %*% t(x)

    return (list(gain = k, parest = h, invcor = P, err_pri = e_pri, err_pos = e_pos, var_pos = v_pos))
}

=====
# common color sets
col_EB_RLS <- rgb (1,133,113, maxColorValue = 255)
col_EB_RLS_line <- rgb (128,205,193, maxColorValue = 255)
col_EB_RLS_lim <- rgb (230,97,1, maxColorValue = 255)
col_EB_MLR <- rgb (166,97,26, maxColorValue = 255)
col_EB_MLR_line <- rgb (223,194,125, maxColorValue = 255)
col_EB_MLR_lim <- rgb (230,97,1, maxColorValue = 255)
col_C <- rgb (49, 163, 84, maxColorValue = 255)
col_H <- rgb (222, 45, 38, maxColorValue = 255)
col_E <- rgb (49, 130, 189, maxColorValue = 255)
col_g1 <- rgb (208,28,139, maxColorValue = 255)
col_g2 <- rgb (51,204,255, maxColorValue = 255)
col_cusum_lim <- rgb (230,97,1, maxColorValue = 255)

dWthreshold <- 0.01
CONF.LEVEL <- 0.95
templabel <- expression(paste('Daily avg. temperature [',degree,'F]'))
lwdverythin <- 0.1
lwdthin <- 0.3
lwdmid <- 0.5
lwdthick <- 0.8
lwdbold <- 1.0

ebl_IP <- expression(paste(italic(E["BL"]), " residuals [Btu/day/", ft^2, "]", sep = ""))
ebl_SI <- expression(paste(italic(E["BL"]), " residuals [MJ/day/", m^2, "]", sep = ""))
energyuse_IP <- expression(paste("Energy use [Btu/day/", ft^2, "]", sep = ""))
energyuse_SI <- expression(paste("Energy use [MJ/day/", m^2, "]", sep = ""))
res_stder_IP <- expression(paste(sigma, " [Btu/day/", ft^2, "]", sep = ""))
res_stder_SI <- expression(paste(sigma, " [MJ/day/", m^2, "]", sep = ""))
Tref_label_IP <- expression(paste(italic(T["ref"]), " [",degree,"F]"))
Tref_label_SI <- expression(paste(italic(T["ref"]), " [",degree,"C]"))

```

```

W_label_IP <- expression(paste(italic(W["oa"]^"+"), " ",
bgroup("[",frac(Btu/day/ft^"2",lb["w"]/lb["da"]),"]"))))
W_label_SI <- expression(paste(italic(W["oa"]^"+"), " ",
bgroup("[",frac(MJ/day/m^"2",kg["w"]/kg["da"]),"]"))))
T_label_IP <- expression(paste(italic(T["oa"]), " ", bgroup("[",frac(Btu/day/ft^"2",degree*F),"]"))))
T_label_SI <- expression(paste(italic(T["oa"]), " ", bgroup("[",frac(MJ/day/m^"2",degree*C),"]"))))
Intercept_IP <- expression(paste("Intercept [Btu/day/", ft^2, "]"))
Intercept_SI <- expression(paste("Intercept [MJ/day/", m^2, "]"))

```

```

calc.yrange <- function(.vectorval, .valtype)
{
  maxval <- max(.vectorval, na.rm = TRUE)
  minval <- min(.vectorval, na.rm = TRUE)
  if (is.infinite(maxval)) {maxval <- 1}
  if (is.infinite(minval)) {minval <- -1}
  if (.valtype == 'EBL') {
    r1 <- max(max(c(abs(maxval), abs(minval))),1)
    r2 <- NA
  } else {
    r1 <- maxval
    r2 <- minval
  }
  return (list(max = r1, min = r2))
}

```

```

legend_energy <- function (.xposition, .yposition) {
  leg.txt <- c('Elec', 'Cool', 'Heat')
  linetypes <- c(1, 1, 1)
  pointtypes <- c(-1, -1, -1)
  legendcols <- c(col_E, col_C, col_H)
  leg.size <- 0.7
  legend(.xposition, .yposition, leg.txt, col = legendcols, lty = 1,
        pch = pointtypes, cex = leg.size, ncol = 3, merge = TRUE, box.lty = 0)
}

```

```

legend_2 <- function (.xposition, .yposition) {
  leg.txt <- c(expression(paste('2', sigma)), expression(paste('3', sigma)))
  linetypes <- c(1, 1)
  pointtypes <- c(-1, -1)
  legendcols <- c(col_EB_MLR_lim, col_EB_MLR_lim)
  leg.size <- 0.7
  legend(.xposition, .yposition, leg.txt, col = legendcols, lty = c(6, 2),
        pch = pointtypes, cex = leg.size, ncol = 2, merge = TRUE, box.lty = 0)
}

```

```

legend_RLSres <- function (.xposition, .yposition) {
  leg.txt <- c("CUSUM alarm")
  linetypes <- c(1)
  pointtypes <- c(-1)
  legendcols <- c(col_EB_RLS_lim)
  leg.size <- 0.7
  legend(.xposition, .yposition, leg.txt, col = legendcols, lty = c(1),
        pch = pointtypes, cex = leg.size, ncol = 2, box.lty = 0)
}

```

```

}

legend_cusumstat <- function (.xposition, .yposition) {
  leg.txt <- c('Upper Cusum', 'Lower Cusum', 'Decision interval')
  linetypes <- c(1, 1, 1)
  pointtypes <- c(-1, -1, -1)
  legendcols <- c(col_g1, col_g2, col_cusum_lim)
  leg.size <- 0.7
  legend(.xposition, .yposition, leg.txt, col = legendcols, lty = 1,
        pch = pointtypes, cex = leg.size, ncol = 3, merge = TRUE)
}

# Graph 1: Energy use
# c(bottom, left, top, right)
plotGraph1 <- function()
{
  par(mar = c(0, 4, 0, 4), mgp = c(1.7, 0.5, 0))
  C <- consData$CHW_ori * 1000000 / consData$Area
  H <- consData$HHW_ori * 1000000 / consData$Area
  E <- consData$ELE_ori * 3412.14 / consData$Area
  Crange <- calc.yrange(C, 'C')$max
  Hrange <- calc.yrange(H, 'H')$max
  Erange <- calc.yrange(E, 'E')$max
  Yrange <- max(Crange, Hrange, Erange)
  Ymin <- 0
  Ymax <- ceiling(Yrange * 1.2/100) * 100
  plot(as.Date(ST, tz=""), 0, type = 'n',
        xlim = c(as.Date(ST), as.Date(ED)),
        ylim = c(Ymin, Ymax),
        xlab = "", ylab = energyuse_IP, axes = F)
  par (tck = 1, lwd = lwdmid)
  axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
  axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
            format="%b%y", labels=FALSE, col = 8, lwd = lwdthin)
  x <- consData$Date
  lines(x, E, type = 'l', col = col_E)
  lines(x, C, type = 'l', col = col_C)
  lines(x, H, type = 'l', col = col_H)
  legend_energy(as.Date(ST), Ymax)
  #par(usr = c(par('usr')[1:2], 0, Ymax * 1055.056 / 0.092903 / 1000000))
  axis(4, at=c(Ymin, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=1),
round(Ymax * 1055.056 / 0.092903 / 1000000, digits=1)), tck= -.02, lwd = lwdmid, las = 1)
  mtext(energyuse_SI, side=4, line=1.2, cex=0.6, las=0)
  box()
}

# Graph 2 --- EBL residuals plot (combined)
plotGraph2 <- function()
{
  EBrIs.err <- EstOut$Pri_Err
  EBresabsrange <- min(calc.yrange(EBrIs.err[(ini.learning + 1):length(EBrIs.err)], 'EBL')$max, 3000)
  par(mar = c(0, 4, 0, 4), mgp = c(1.7, 0.5, 0))
  plot(as.Date(ST, tz=""), 0, type = 'n',

```

```

xlim = c(as.Date(ST), as.Date(ED)),
ylim = c(-1 * ceiling(EBresabsrange * 1.2/100) * 100, ceiling(EBresabsrange *
1.2/100) * 100),
xlab = "", ylab = ebl_IP, axes = F)
par (tck = 1, lwd = lwdmid)
axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

# MLR results
if (!is.na(lm_out[[1]])) {
  # plot original data - predicted for the screening month (the last month)
  # adjust the length of cons_pred and lm_out$outdata in case there are missing data in the
beginning of the period.
  if (nrow(cons_pred) < nrow(lm_out$outdata)) {
    pST <- cons_pred$Date[1]
    pED <- cons_pred$Date[nrow(cons_pred)]
    modoutdata <- selectTime(lm_out$outdata, pST, pED)
    y <- cons_pred$EBL_ori - modoutdata$pred.fit
    x <- cons_pred$Date
  } else {
    x <- lm_out$outdata$Date
    y <- cons_pred$EBL_ori - lm_out$outdata$pred.fit
  }
  lines(x, y, type = 'l', col = col_EB_MLR_line)
  points(x, y, pch=3, col = col_EB_MLR, cex = 0.2)
  # plot approximate prediction intervals (2-sigma) for MLR model based on one year data
  curve(x * 0 + lm_out$RMSE * 2, add = T, lwd = lwdmid, col = col_EB_MLR_lim)
  curve(x * 0 - lm_out$RMSE * 2, add = T, lwd = lwdmid, col = col_EB_MLR_lim)
  curve(x * 0 + lm_out$RMSE * 3, add = T, lwd = lwdmid, lty = 2, col =
col_EB_MLR_lim)
  curve(x * 0 - lm_out$RMSE * 3, add = T, lwd = lwdmid, lty = 2, col =
col_EB_MLR_lim)
}

# RLS results
x <- EstOut$Date
lines(x, EBrls.err, type = 'l', col = col_EB_RLS_line)
points(x, EBrls.err, pch=3, col = col_EB_RLS, cex = 0.2)

alarmedDays <- EstOut[EstOut$salarm == 1, ]
if (nrow(alarmedDays) > 0){
  for (n in 1:nrow(alarmedDays)) {
    abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col = 'blue')
  }
}

legend_2(as.Date(ST) + (as.Date(ED) - as.Date(ST))/2, ceiling(EBresabsrange * 1.2/100) * 100)

abline(h = 0)
box(col = 9)
}

```

```

# EBL residuals plot (MLR)
plotGraph_MLR_res <- function()
{
  EBrIs.err <- EstOut$Pri_Err
  EBresabsrange <- min(calc.yrange(EBrIs.err[(ini.learning + 1):length(EBrIs.err)], 'EBL')$max, 3000)
  par(mar = c(0, 4, 0, 4), mgp = c(1.7, 0.5, 0))
  Ymin <- -1 * ceiling(EBresabsrange * 1.2/100) * 100
  Ymax <- ceiling(EBresabsrange * 1.2/100) * 100
  plot(as.Date(ST, tz=""), 0, type = 'n',
        xlim = c(as.Date(ST), as.Date(ED)),
        ylim = c(Ymin, Ymax),
        xlab = "", ylab = ebl_IP, axes = F)
  par (tck = 1, lwd = lwdmid)
  axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
  axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
            format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

  # MLR results
  if (!is.na(lm_out[[1]])) {
    # plot original data - predicted for the screening month (the last month)
    # adjust the length of cons_pred and lm_out$outdata in case there are missing data in the
    # beginning of the period.
    if (nrow(cons_pred) < nrow(lm_out$outdata)) {
      pST <- cons_pred$Date[1]
      pED <- cons_pred$Date[nrow(cons_pred)]
      modoutdata <- selectTime(lm_out$outdata, pST, pED)
      y <- cons_pred$EBL_ori - modoutdata$pred.fit
      x <- cons_pred$Date
    } else {
      x <- lm_out$outdata$Date
      y <- cons_pred$EBL_ori - lm_out$outdata$pred.fit
    }
    lines(x, y, type = 'l', col = col_EB_MLR_line)
    points(x, y, pch=3, col = col_EB_MLR, cex = 0.2)
    # plot approximate prediction intervals (2-sigma) for MLR model based on one year data
    x <- cons_pred$Date
    curve(x * 0 + lm_out$RMSE * 2, add = T, lwd = lwdmid, lty = 6, col =
col_EB_MLR_lim)
    curve(x * 0 - lm_out$RMSE * 2, add = T, lwd = lwdmid, lty = 6, col =
col_EB_MLR_lim)
    curve(x * 0 + lm_out$RMSE * 3, add = T, lwd = lwdmid, lty = 2, col =
col_EB_MLR_lim)
    curve(x * 0 - lm_out$RMSE * 3, add = T, lwd = lwdmid, lty = 2, col =
col_EB_MLR_lim)
  }
  legend_2(as.Date(ED)-112, Ymax)
  axis(4, at=c(Ymin, 0, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=1),
0, round(Ymax * 1055.056 / 0.092903 / 1000000, digits=1)), tck = -.02, lwd = lwdmid, las = 1)
  mtext(ebl_SI, side=4, line=1.2, cex=0.6, las=0)
  abline(h = 0)
  box(col = 9)
}

```

```

# EBL residuals plot (RLS)
plotGraph_RLS_res <- function()
{
  EBrls.err <- EstOut$Pri_Err
  EBresabsrange <- min(calc.yrange(EBrls.err[(ini.learning + 1):length(EBrls.err)], 'EBL')$max, 3000)
  Ymin <- -1 * ceiling(EBresabsrange * 1.2/100) * 100
  Ymax <- ceiling(EBresabsrange * 1.2/100) * 100
  par(mar = c(0, 4, 0, 4), mgp = c(1.7, 0.5, 0))
  plot(as.Date(ST, tz=""), 0, type = 'n',
       xlim = c(as.Date(ST), as.Date(ED)),
       ylim = c(Ymin, Ymax),
       xlab = "", ylab = ebl_IP, axes = F)
  par (tck = 1, lwd = lwdmid)
  axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
  axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
            format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

  # RLS results
  x <- EstOut$Date
  lines(x, EBrls.err, type = 'l', col = col_EB_RLS_line)
  points(x, EBrls.err, pch=3, col = col_EB_RLS, cex = 0.2)

  alarmedDays <- EstOut[EstOut$alarm == 1, ]
  if (nrow(alarmedDays) > 0){
    for (n in 1:nrow(alarmedDays)) {
      abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col = col_EB_RLS_lim)
    }
  }
  legend_RLSres(as.Date(ED)-112, Ymax)
  axis(4, at=c(Ymin, 0, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=1),
0, round(Ymax * 1055.056 / 0.092903 / 1000000, digits=1)), tck= -.02, lwd = lwdmid, las = 1)
  mtext(ebl_SI, side=4, line=1.2, cex=0.6, las=0)
  abline(h = 0)
  box(col = 9)
}

# Graph TestStat
plotGraph_TestStat <- function()
{
  par(mar = c(0, 4, 0, 4), mgp = c(1.7, 0.5, 0))
  plot(as.Date(ST, tz=""), 0, type = 'n',
       xlim = c(as.Date(ST), as.Date(ED)),
       ylim = c(ceiling(h * 1.2/10) * 10 * -1, ceiling(h * 1.2/10) * 10),
       xlab = "", ylab = 'CUSUM statistics', axes = F)
  par (tck = 1, lwd = lwdmid)
  axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
  axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
            format="%b%y", labels = TRUE, col = 8, lwd = lwdthin)

  # g_1
  x <- EstOut$Date
  y <- EstOut$CUSUM_g_1

```

```

lines(x, y, type = 'l', col = col_g1)

# g_2
y <- EstOut$CUSUM_g_2 * -1
lines(x, y, type = 'l', col = col_g2)

# limit line
curve(x * 0 + h, add = T, lwd = lwdmid, col = col_cusum_lim)
curve(x * 0 - h, add = T, lwd = lwdmid, col = col_cusum_lim)
#legend_cusumstat(as.Date(ST) + (as.Date(ED) - as.Date(ST))/2, ceiling(h * 1.2/10) * 10)
text(as.Date(ST), h, "h = 5", pos = 3)
text(as.Date(ST), -1 * h, "h = 5", pos = 1)
abline(h = 0)
box(col = 9)

#arrows
pos_arrows <- as.Date(ST) + 28
arrows(c(pos_arrows, pos_arrows), c(0, 0), c(pos_arrows, pos_arrows), c(9, -9), angle = 30,
length = 0.04, code = 2)
upperlabel <- expression(italic(C^"+"))
lowerlabel <- expression(italic(-C^-"))
text(pos_arrows, 2.5, upperlabel, pos=2)
text(pos_arrows, -2.5, lowerlabel, pos=2)

}

# Graph_Var --- Variance
plotGraph_Var <- function()
{
y <- sqrt(EstOut$EmpiricalVar)
Ymin <- 0
Ymax <- max(lm_out$RMSE * 2, calc.yrange(y, 'var')$max)
par(mar = c(0, 5, 0, 5), mgp = c(1.7, 0.5, 0))
plot(as.Date(ST, tz=""), 0, type = 'n',
      xlim = c(as.Date(ST), as.Date(ED)),
      ylim = c(Ymin, Ymax),
      xlab = "", ylab = res_stderr_IP, axes = F)
par (tck = 1, lwd = lwdmid)
axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
          format="%b%y", labels = TRUE, col = 8, lwd = lwdthin)

# EWMA smoothing
x <- EstOut$Date
lines(x, y, type = 'l', col = col_EB_RLS)

# MLR sigma
x <- lm_out$outdata$Date
y <- rep(1, length(x)) * lm_out$RMSE
lines(x, y, type = 'l', col = col_EB_MLR)

# alarm
alarmedDays <- EstOut[EstOut$alarm == 1, ]

```

```

    if (nrow(alarmedDays) > 0){
      for (n in 1:nrow(alarmedDays)) {
        abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col = col_EB_RLS_lim)
      }
    }

    #SI labels
    axis(4, at=c(Ymin, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=2),
round(Ymax * 1055.056 / 0.092903 / 1000000, digits=2)), tck= -.02, lwd = lwdmid, las = 1)
    mtext(res_stderr_SI, side=4, line=1.5, cex=0.6, las=0)
    box(col = 9)
  }

# Graph 3 --- Tref
plotGraph_Tref <- function()
{
  y <- EstOut$T_ref
  Ymin <- 0
  Ymax <- 120
  par(mar = c(0, 5, 0, 5), mgp = c(1.7, 0.5, 0))
  plot(as.Date(ST, tz=""), 0, type = 'n',
        xlim = c(as.Date(ST), as.Date(ED)),
        ylim = c(Ymin, Ymax),
        xlab = "", ylab = Tref_label_IP, axes = F)
  par (tck = 1, lwd = lwdmid)
  axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
  axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
            format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

  # EBL estimates
  if (!is.na(lm_out[[1]])) {
    x <- EstOut$Date
    curve(x * 0 + (lm_out$Est_int / lm_out$Est_T * -1), add = T, lwd = lwdmid, col =
col_EB_MLR)
  }

  # RLS results
  x <- EstOut$Date
  lines(x, y, type = 'l', col = col_EB_RLS)
  # alarm
  alarmedDays <- EstOut[EstOut$alarm == 1, ]
  if (nrow(alarmedDays) > 0){
    for (n in 1:nrow(alarmedDays)) {
      abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col =
col_EB_RLS_lim)
    }
  }

  axis(4, at=c(Ymin, Ymax), labels=c(round((Ymin-32) * 5 / 9, digits=1), round((Ymax-32) * 5 / 9,
digits=1)), tck= -.02, lwd = lwdmid, las = 1)
  mtext(Tref_label_SI, side=4, line=1.5, cex=0.6, las=0)

  box(col = 9)

```



```

}

# Intercept
plotGraph_Int <- function()
{
y <- EstOut$Est_Int
par(mar = c(0, 5, 0, 5), mgp = c(1.7, 0.5, 0))
Ymin <- min(0, calc.yrange(y[(ini.learning + 1):length(y)], 'Int')$min)
Ymax <- calc.yrange(y[(ini.learning + 1):length(y)], 'Int')$max
plot(as.Date(ST, tz=""), 0, type = 'n',
      xlim = c(as.Date(ST), as.Date(ED)),
      ylim = c(Ymin, Ymax),
      xlab = "", ylab = Intercept_IP, axes = F)
par (tck = 1, lwd = lwdmid)
axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
          format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

# RLS estimates
x <- EstOut$Date
lines(x, y, type = 'l', col = col_EB_RLS)

# EBL estimates
if (!is.na(lm_out[[1]])) {
  curve(x * 0 + lm_out$Est_int, add = T, lwd = lwdmid, col = col_EB_MLR)
}

# alarm
alarmedDays <- EstOut[EstOut$alarm == 1, ]
if (nrow(alarmedDays) > 0){
  for (n in 1:nrow(alarmedDays)) {
    abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col =
col_EB_RLS_lim)
  }
}

axis(4, at=c(Ymin, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=1),
round(Ymax * 1055.056 / 0.092903 / 1000000, digits=1)), tck = -.02, lwd = lwdmid, las = 1)
mtext(Intercept_SI, side=4, line=1.5, cex=0.6, las=0)
box(col = 9)
}

# Temperature slope
plotGraph_Tslope <- function()
{
y <- EstOut$Est_T
par(mar = c(0, 5, 0, 5), mgp = c(2, 0.3, 0))
Ymin <- calc.yrange(y[(ini.learning + 1):length(y)], 'Tref')$min
Ymax <- max(0, calc.yrange(y[(ini.learning + 1):length(y)], 'Tref')$max)
plot(as.Date(ST, tz=""), 0, type = 'n',
      xlim = c(as.Date(ST), as.Date(ED)),
      ylim = c(Ymin, Ymax),

```

```

        xlab = ", ylab = T_label_IP, axes = F)
par (tck = 1, lwd = lwdmid)
axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
          format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

# RLS estimates
x <- EstOut$Date
lines(x, y, type = 'l', col = col_EB_RLS)

# EBL estimates
if (!is.na(lm_out[[1]])) {
  curve(x * 0 + lm_out$Est_T, add = T, lwd = lwdmid, col = col_EB_MLR)
}

# alarm
alarmedDays <- EstOut[EstOut$salarm == 1, ]
if (nrow(alarmedDays) > 0){
  for (n in 1:nrow(alarmedDays)) {
    abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col =
col_EB_RLS_lim)
  }
}

axis(4, at=c(Ymin, Ymax), labels=c(round(Ymin * 0.020441759, digits=4), round(Ymax *
0.020441759, digits=5)), tck= -.02, lwd = lwdmid, las = 1)
mtext(T_label_SI, side=4, line=4, cex=0.6, las=0)
box(col = 9)
}
# Humidity slope
plotGraph_W <- function()
{
y <- EstOut$Est_dW
par(mar = c(0, 5, 0, 5), mgp = c(2, 0.3, 0))
Ymin <- calc.yrange(y[(ini.learning + 1):length(y)], 'dW')$min
Ymax <- max(0,calc.yrange(y[(ini.learning + 1):length(y)], 'dW')$max)
plot(as.Date(ST, tz=""), 0, type = 'n',
      xlim = c(as.Date(ST), as.Date(ED)),
      ylim = c(Ymin, Ymax),
      xlab = ", ylab = W_label_IP, axes = F)
par (tck = 1, lwd = lwdmid)
axis(2, lty = 1, col = 8, lwd = lwdthin, las = 1)
axis.Date(1, at=seq.Date(as.Date(ST),as.Date(ED), by="1 month"),
          format="%b%y", labels = FALSE, col = 8, lwd = lwdthin)

# RLS estimates
x <- EstOut$Date
lines(x, y, type = 'l', col = col_EB_RLS)

# EBL estimates
if (!is.na(lm_out[[1]])) {
  curve(x * 0 + lm_out$Est_W, add = T, lwd = lwdmid, col = col_EB_MLR)
}

```

```

# alarm
alarmedDays <- EstOut[EstOut$alarm == 1, ]
if (nrow(alarmedDays) > 0){
  for (n in 1:nrow(alarmedDays)) {
    abline(v = alarmedDays$Date[[n]], lwd = lwdthin, lty = 1, col = col_EB_RLS_lim)
  }
  axis(4, at=c(Ymin, Ymax), labels=c(round(Ymin * 1055.056 / 0.092903 / 1000000, digits=1),
round(Ymax * 1055.056 / 0.092903 / 1000000, digits=1)), tck= -.02, lwd = lwdmid, las = 1)
  mtext(W_label_SI, side=4, line=4, cex=0.6, las=0)
  box(col = 9)
}

```