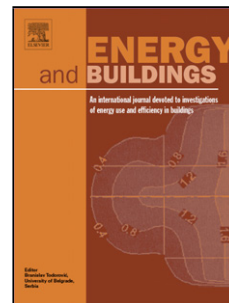


Accepted Manuscript

Title: Estimation models of heating energy consumption in schools for local authorities planning

Author: Alfonso Capozzoli Daniele Grassi Francesco Causone



PII: S0378-7788(15)30136-5
DOI: <http://dx.doi.org/doi:10.1016/j.enbuild.2015.07.024>
Reference: ENB 6017

To appear in: *ENB*

Received date: 3-3-2015
Revised date: 2-7-2015
Accepted date: 11-7-2015

Please cite this article as: A. Capozzoli, D. Grassi, F. Causone, Estimation models of heating energy consumption in schools for local authorities planning, *Energy and Buildings* (2015), <http://dx.doi.org/10.1016/j.enbuild.2015.07.024>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

- The annual heating energy consumptions of eighty school buildings are analysed
- Two energy estimation models were developed to support public authorities planning
- A multiple regression model was built using nine different influencing variables
- CART enables also non-expert users to extract information for decision making
- MAE, RMSE and MAPE were calculated to compare the performance of estimation models

Accepted Manuscript

1 Estimation models of heating energy consumption in schools 2 for local authorities planning

3
4 Alfonso Capozzoli^{1*}, Daniele Grassi¹, Francesco Causone²

5
6 1. TEBE Research group, Department of Energy, Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129 Torino,
7 Italy.

8 2. Department of Energy, Politecnico di Milano, Via Lambruschini 4, 20156 Milano, Italy

9 * Corresponding author: Alfonso Capozzoli

10 Tel.: +39 011 0904413; fax: +39 011 0904499

11 E-mail addresses: alfonso.capozzoli@polito.it (A. Capozzoli), grassi.daniele88@gmail.com (D. Grassi),

12 francesco.causone@polimi.it (F. Causone)

13 Keywords: School buildings; Multiple Linear Regression; Heating energy estimation; Decision Tree; Public authorities
14 planning

16 Abstract

17 Large building stocks should be well managed, in terms of ordinary activities and formulating strategic plans, to achieve
18 energy savings through increased efficiency. It is becoming extremely important to have the capability to quickly and
19 reliably estimate buildings' energy consumption, especially for public authorities and institutions that own and manage
20 large building stocks. This paper analyses the heating energy consumption of eighty school buildings located in the
21 north of Italy. Two estimation models are developed and compared to assess energy consumption: a Multiple Linear
22 Regression (MLR) model and a Classification and Regression Tree (CART). The CART includes interpretable decision
23 rules that enable non-expert users to quickly extract useful information to benefit their decision making. The output of
24 MLR model is an equation that accounts for all of the major variables affecting heating energy consumption. Both
25 models were compared in terms of Mean Absolute Error (MAE), Root Mean Square error (RMSE), and Mean Absolute
26 Percentage error (MAPE). The analysis determined that the heating energy consumption of the considered school
27 buildings was mostly influenced by the gross heated volume, heat transfer surfaces, boiler size, and thermal
28 transmittance of windows.

29

30 Nomenclature

31 β	Estimated Coefficient of Multiple Linear Regression Model
32 CART	Classification and Regression Tree
33 D-W	Durbin-Watson test
34 E	Error associated to the tree
35 EUI	Energy Use Intensity
36 EUI_{st}	Standard Energy Use Intensity

37	$EUI_{st,s}$	Standard and Specific Energy Use Intensity
38	F	Fisher-Snedecor test
39	HDD_{con}	Conventional Heating Degree Days
40	HDD_{real}	Real Heating Degree Days
41	MAE	Mean Absolute Error
42	MAPE	Mean Absolute Percentage Error
43	MLR	Multiple Linear Regression
44	POW	Boiler Size (Heat Input)
45	R^2	Coefficient of Determination
46	R^2_{adj}	Adjusted Coefficient of Determination
47	RMSE	Root Mean Square Error
48	SUR	Heat Transfer Surface
49	$U_{windows}$	Thermal Transmittance of Windows
50	VIF	Variance Inflation Factors
51	VOL	Heated Gross Volume
52		

53 1. Introduction

54 1.1 Energy consumption analysis in school buildings

55 Buildings are responsible for about 40 % of the total energy consumption in developed countries [1]. In countries like
 56 Italy, about 60 % of the existing building stock is more than 40-years-old [2]. A rapid and substantial energy retrofit
 57 program is therefore required for these existing buildings. There are about two million households in Italy that live in
 58 buildings requiring either demolition and rebuilding or refurbishment. Directive 2010/31/EU (EPBD recast) requires
 59 buildings, or parts of buildings, to meet a minimum energy performance or be subject to a retrofit or refurbishment.
 60 These requirements may be met by renovating a building's envelope and systems, but an effective management can also
 61 significantly impact a building's energy consumption. The EU Directive [3, 4, 5] requires public buildings to play an
 62 exemplary role in terms of energy savings. Public utilities and buildings that are typically owned and managed by
 63 municipalities include: street lighting, schools, administrative buildings, public transport, and sport centres, such as
 64 swimming pools and gymnasiums [6]. Local governments would clearly benefit from having access to energy
 65 consumption data. Further, being able to understand the savings potential of these assets would help to prioritise energy
 66 and environmental projects and better illuminate their financial aspects [7]. According to the US Department of Energy,
 67 school buildings constitute a major part of the public building stock. Around 25 % of the energy expenses in schools
 68 could be saved through better building designs and more energy-efficient technologies, combined with improvements in
 69 operation and maintenance [8].

70 De Santoli et al. [9] evaluated the energy performance of public schools in Rome. They defined intervention strategies
 71 to reduce energy consumption and identified action priorities by means of a simple payback time analysis (PBT).
 72 Dimoudi et al. [10] conducted an energy simulation to study the energy savings potential of school buildings in Greece.
 73 Kim Tae Woo et al. [11] analysed the energy consumption of some elementary schools in South Korea by utilising
 74 monitoring data from January 2006 to December 2010. They determined that electrical energy was consumed the most,
 75 followed by gas and oil. During the monitoring period, electrical energy continued to increase its relevance on the
 76 energy breakdown because of cooling/heating system replacements. These and other studies were carried out in recent
 77 years to estimate the energy consumption of school buildings. The literature shows that there are two main approaches
 78 for estimating a building's energy consumption: the direct approach or the inverse approach. The first approach
 79 calculates the energy demand by running an energy simulation under a steady state or dynamic conditions. The second
 80 approach uses historical data to produce data driven models that estimate the energy consumption.

81 A large part of the current literature focuses on the inverse approach. Analysts and decision makers have access to
82 several applications of new or recast versions of existing models. Corrado et al. [12] defined a simplified method for
83 predicting future consumption based on climatic and real use data on a stock of 120 school buildings. Corgnati et al.
84 [13] then validated this method, using another stock of 118 schools, as did Ariaudio et al. [14]. Amber et al. [15]
85 gathered daily values of a school building's electrical consumption on the Southwark campus of the London South
86 Bank University from 2007 to 2013 and then developed a multiple regression model to estimate future daily electrical
87 consumption. Beusker et al. [16] evaluated the energy consumption of schools and sports facilities in Germany using
88 different linear and nonlinear regression models. Thewes et al. [17] presented a regression model with categorical
89 variables to predict the electrical and heating energy consumption of school buildings in Luxembourg.
90 Innovative techniques, including machine learning, data mining, and knowledge discovery in databases, have also been
91 successfully applied to building energy consumption data in recent years [18]. In particular, a classification tree which
92 consists of a multi-stage decision-making process that is useful to categorise observations in a finite number of classes,
93 can be a powerful estimation tool. This method has not yet been applied in other studies to estimate the energy
94 consumption of school buildings.
95 In this paper, the heating energy consumption of a school building stock located in the north of Italy is analysed using a
96 Multiple Linear Regression (MLR) model and a Classification and Regression Tree (CART). Both MLR model and
97 CART are data driven models that have been successfully applied to estimate a building's energy demand.
98 Nevertheless, the outcome of MLR is an equation, while the output of CART are decision rules that allow users to
99 quickly extract relevant information [19]. This characteristic substantially changes the practical applicability of the two
100 models.

101 **1.2 Implementation of multiple regression analysis and classification tree for buildings' energy use estimations**

102 In recent years, numerous researchers successfully employed multiple regression model as a tool for energy
103 consumption estimations. Al-Garni et al. [20] correlated electrical energy consumption with relevant climatic variables
104 (air temperature, relative humidity, solar radiation), and variable occupant populations through statistical methods
105 (regression model) to forecast the overall electrical energy consumption in Eastern Saudi Arabia. Aranda et al. [21]
106 developed three regression models to predict the Spanish banking sector's annual energy consumption. The first model
107 can be used to estimate the energy consumption of the whole banking sector, while the second estimates the energy
108 consumption for branches under conditions of a low severity winter climate and the third under conditions of a high
109 severity winter climate. The variance reported for the three models is 58 %, and 68 %, respectively. Korolija et al. [22]
110 developed regression models to predict the annual heating, cooling, and electrical auxiliary energy consumption of five
111 different types of HVAC systems (variable air volume – VAV, constant air volume – CAV, fan-coil system with
112 dedicated air (FC), and two chilled ceiling systems with dedicated air, radiator heating, and either embedded pipes –
113 EMB - or exposed aluminium panels – ALU) for office buildings in the UK. Freire et al. [23] used independent
114 variables like energy consumption, ventilation and air conditioning power, outdoor temperature, relative humidity, and
115 total solar radiation to develop a regression equation to predict the indoor air temperature and relative humidity for two
116 buildings with low and high thermal mass. The literature demonstrates therefore that regression models offer a robust
117 methodology for estimating a building's energy consumption (e.g., heating, cooling, lighting, etc.).
118 Decision trees belong to the *machine learning algorithms* family. This method is recognised as an emerging analysis
119 tool and is currently receiving plenty of attention from applied research. Yu et al. [18] used the decision tree to classify
120 and predict building energy consumption. This method was applied to Japanese residential buildings for predicting and

121 classifying building Energy Use Intensity (EUI) levels based on training data. This tool was then evaluated on a sample
122 test.

123 Zhao et al. [24] used a C4.5 decision tree algorithm, locally weighted naïve Bayes and support vector machine, to
124 classify occupant behaviour and to create schedule models for building energy simulation. The results show that the
125 C4.5 algorithm correctly classified 90 % of individual behaviour and this allowed getting closer to the real group
126 schedule. Mikučionienė et al. [25] used a decision tree to increase the sustainability and improve the criteria for
127 evaluating energy efficiency measures in a public building renovation in Lithuania. By analysing and weighting each
128 variable (related to insulation of external walls, roof insulation, heating substation renovation, reconstruction of the
129 entire heating system, and installation of a ventilation system with exhaust air heat recovery), the researchers created a
130 decision tree to evaluate the influence of each variable on energy consumption. The results show that this algorithm
131 reduces the amount of data that must be understood by transforming it into a more compact form while still preserving
132 the basic substance. The researchers determine whether the data are characterised by well-separated object classes and
133 finally, this algorithm determines the precise relationship between attributes and their class.

134 In this paper, two different estimation models are developed using a database consisting of 80 school buildings located
135 in the province of Turin. The estimation models include climatic, envelope and heating system variables, and annual
136 metered heating energy consumption. They are:

- 137 - a Multiple Linear Regression (MLR) model that estimates the energy use for heating based on geometrical,
138 climate, and thermo-physical characteristics. This model creates an equation that relates n independent
139 variables to the dependent variable;
- 140 - a Classification and Regression Tree (CART) which consists of a multi-stage decision-making process to
141 classify observations in a finite number of classes. The model's output is a flowchart constructed by
142 subdividing the observations into homogeneous subsets with respect to the dependent variable or response
143 (represented in our model by heating energy consumption).

144 The two estimation models are compared to determine which one is more accurate in terms of a residuals analysis and
145 errors (MAE, RMSE, MAPE). The possibilities and limitations of the two models are ultimately contrasted,
146 highlighting advantages and disadvantages for their use by a final operator, such as a consultant or a decision maker.
147 Moreover, this paper discusses the practical application and robustness of the constructed estimation models.

148 **2 Methodology**

149 A wide range of theoretical and practical factors that are relevant to each building should be considered to create
150 estimation models that analyse building energy consumption. The methodology followed in this study is schematised as
151 shown in the flowchart in Figure 1.

152 An existing database was initially analysed to evaluate the consistency of the school building stock. The available
153 variables are associated with a building's envelope, heating/cooling systems, and location. This step is useful to
154 understand the limits of applicability of the models, which may be applied to other building stocks with similar features,
155 once they are validated. In the second step of the analysis, two estimation models were implemented (MLR model and
156 CART). Finally, in the third phase, the two models were developed and compared, highlighting their usefulness for
157 public school managers.

158 2.1 Pre-processing analysis

159 The database contains information from 80 school buildings without sport facilities (sport halls), situated in the
160 Province of Turin (Italy). The initial dataset was composed of 120 school buildings located in the same area, but the
161 sample was reduced to 80 schools due to missing heating energy consumption data from 40 school buildings. The
162 analysed influencing variables are related to the opaque and transparent building envelope, heating systems, building
163 geometry features, and climatic data.

164 From a climatic point of view, the Province of Turin is located in the Italian climate zones E and F. The analysed
165 buildings are located in a climate with Conventional Heating Degree Days (HDD_{conv}) ranging from 2517 to 3197 DD.
166 Figure 2 shows the frequency distribution for the gross heated volume and heat transfer surface of the sampled
167 buildings. The majority of schools have a gross heated volume lower than 35000 m^3 (about 60 %). Schools with a
168 higher gross heated volume are composed of two or more buildings. In addition, about 60 % of the sampled schools
169 have values of heat transfer surface lower than 10000 m^2 . For this reason, most of the sample is composed of buildings
170 with an aspect ratio (ratio of heat transfer surface on gross heated volume) range from 0.25 to 0.40 m^{-1} . The heat losses
171 mainly depend on the quality of the building envelope and not from the building shape.

172 Figure 3 shows the frequency distribution of the sampled buildings for the thermal transmittance of walls and windows.
173 As can be seen, most of the buildings are characterised by a thermal transmittance of windows higher than $4 \text{ W}/(\text{m}^2 \text{ K})$
174 (about 65 % of the sample is composed of single glazing) and by opaque walls without thermal insulation (80 % of the
175 sample is characterised by values higher than $0.40 \text{ W}/\text{m}^2\text{K}$).

176 Figure 4 shows the frequency distribution of the sampled buildings for boiler size (heat input) and average system
177 efficiency. The boiler size (heat input) ranges from values lower than 500 kW to values higher than 8000 kW. 12 % of
178 the schools are equipped with a boiler size lower than 500 kW, 43 % from 500 to 1500 kW, 17 % from 1500 to 2000
179 kW, and only 28 % by a boiler size higher than 2000 kW. Analysing the frequency distribution of the average seasonal
180 system efficiency (the ratio between building energy need and primary energy) reveals that 83 % of the sample have
181 values lower than 0.70. This figure denotes the presence of high thermal losses in subsystems. Moreover, the schools
182 are equipped with old emission subsystems (cast iron radiators), old distribution subsystems (non-insulated pipes), and
183 old control subsystems, i.e. centralised control that is only installed at the generation system level (e.g. climatic control).

184 Several variables related to school buildings should be considered in a comprehensive analysis of energy consumption.
185 Ventilation rates, hours of use, set points and time clock settings, infiltration rates, internal heat gains, solar gains,
186 geometrical building characteristics, building envelope physical variables, heating system features, outdoor temperature,
187 and number of pupils and classes are all considered important variables in characterising a school's energy use.
188 Moreover, occupant behaviour can significantly impact energy consumption, particularly the opening and closing of
189 windows. However, some of these variables (e.g. infiltration and ventilation rates or variables related to occupant
190 behaviour) are very difficult to obtain.

191 In [26], it was claimed that the floor surface and/or the volume (mostly the volume) primarily influenced the heating
192 energy consumption and the electrical energy in school buildings in their analysed sample [27]. In some cases, it was
193 also verified that the data related to the transmittance of opaque components of the façades, the boiler size, and the daily
194 period of use significantly influenced the heating energy consumption [27].

195 In our work the available data collected for the analysed sample to characterise the heating energy consumption for each
196 school building are: *real heating degree day*, *gross heated volume*, *heat transfer surface*, *aspect ratio*, *floor heated*

197 *area, building height, numbers of floors, thermal transmittance of walls, thermal transmittance of windows, boiler size*
 198 *(heat input), number of classrooms, number of pupils, annual operating time, average seasonal system efficiency.*

199 Table 1 provides a list of the variables with the definition of the data location, central tendency, and dispersion for each
 200 of them. From the literature [26, 27], we know that the selected variables can be considered as the most influential
 201 factors. No major effect of controls may be reported in the database, since no local control was installed in the school.
 202 Occupant behaviour definitely affected the final energy performance by opening and closing windows. However,
 203 information about occupant behaviour were not available and are very difficult to get.

204 In order to standardise the impact of climate on heating energy consumptions, the degree-days method was applied [28].
 205 For this purpose, standard heating energy consumptions (EUI_{st} - Standard Energy Use Intensity) for each school
 206 buildings was defined as:

$$207 \quad EUI_{st} = EUI \cdot (HDD_{conv} / HDD_{real}) \quad (1)$$

208 where EUI is the heating energy consumption [kWh], HDD_{con} are the Heating Conventional Degree Days, and HDD_{real}
 209 are the Heating Real Degree Days (related to the year 2012). The average heating energy consumption of the school
 210 building stock EUI_{st} is equal to 830 MWh/year (Figure 5). In order to compare buildings of different sizes, the EUI_{st}
 211 was normalised by the gross heated volume ($EUI_{st,s}$ - Standard and Specific Energy Use Intensity). The $EUI_{st,s}$ ranges
 212 from 17.74 to 61.12 kWh/m³/year (Figure 5).

213 2.2 Outliers detection

214 A pre-processing analysis [29] is required to identify outliers before creating an estimation model. An observation is an
 215 outlier when it departs from other members of the sample and appears to be inconsistent with the remaining dataset. The
 216 presence of one or more outliers could reduce the capacity of the models to estimate the heating energy consumption.
 217 Outliers should be eliminated from the dataset [30, 31], however, their treatment is not simple. Several indexes were
 218 evaluated to identify outliers in the dataset. In fact, these indexes can be used together to perform an accurate screening
 219 of the database:

- 220 - z-score;
- 221 - Mahalanobis Distance;
- 222 - Index of Mardia;
- 223 - Distances Cook;
- 224 - Leverage Value.

225 These techniques made it possible to detect outliers at both multivariate and univariate levels. The first index (z-score)
 226 detected outliers at the univariate level, i.e. for each variable. The other four indexes detected multivariate outliers by
 227 considering a combination of different variables. These outlier detection categories are complementary and should be
 228 used together. Indeed, a case cannot be considered an outlier if it only has one single distorted value. At the same time,
 229 multivariate outliers represent a pattern of responses that are unlikely to be comparable to the rest of the sample.

230 The definitions of the analysed indexes are briefly explained in the following.

231 The z-score is used to measure when the observed value deviates from the mean value. It is expressed by the mean of
 232 the following equation.

$$233 \quad z - score = (x - \bar{x}) / DS \quad (2)$$

234 where x is the observed value, \bar{x} is the average value, and DS is the standard deviation. On the basis of Chebyshev's
235 theorem, if $z\text{-score} \geq$ the values are potential outliers.

236 The Mahalanobis distance (D_{hk}) is a statistical measure of the distance between the units. It is calculated by taking into
237 account the correlation between variables:

$$238 \quad D_{hk} = \sqrt{(x_h - x_k)^T \cdot W^{-1} \cdot (x_h - x_k)} \quad \text{with } h \neq k = 1, \dots, n \quad (3)$$

239 where x_h and x_k are the vectors with the observations on the samples h and k , and W is the variance-covariance matrix
240 between the observed variables. As a rule of thumb, values of D_{hk} higher than the chi-square critical ($\alpha = 0,001$, degree
241 of freedom = predictors) are considered abnormal points.

242 The Index of Mardia checks if the relationship between variables can be considered linear. The multivariate normality is
243 met if the Index of Mardia is less than a critical value:

$$244 \quad Mah_{critical} = k \cdot (k + 2) \quad (4)$$

245 where k is the number of predictors.

246 The Distances Cook (D_i) is the distance between the regression line that includes all observations and the regression line
247 that does not include the i -th observation:

$$248 \quad D_i = |\hat{y} - \hat{y}_i| / p \cdot DS \quad (5)$$

249 where \hat{y} is the expected value, \hat{y}_i is the expected value without the use of the i -th case, p is the order of the multiple
250 regression analysis, and DS is the standard deviation. Generally, values higher than 1 are considered abnormal points.

251 The Leverage Value ($L_{average}$) is a measure of how much the specified value of the independent variable deviates from
252 its mean. The values vary between zero (no influence) and $(n/(n-1))$ (greatest influence). The average value corresponds
253 to:

$$254 \quad L_{average} = (k+1)/n \quad (6)$$

255 where k is the number of predictors and n is the number of cases analysed. As a rule of thumb, values higher than two
256 or three times the average values are considered abnormal points.

257 **2.3 Multiple Linear Regression model**

258 The multivariate statistical analysis [32, 33] can estimate the value of some variables, if the parameters included in the
259 model are actually relevant for the building's final energy consumption. The MLR model (classical model for parameter
260 estimation) is expressed as follows:

$$261 \quad Y = \beta_0 + \beta_1 \cdot x_1 + \beta_2 \cdot x_2 + \dots + \beta_n \cdot x_n + \varepsilon \quad (7)$$

262 where Y is the dependent variable, β_0 is the intercepts, $\beta_{1,\dots,p}$ are the estimated coefficient of MLR model, and ε is the
263 statistical error. The regression model's coefficients $\beta_{1,\dots,p}$ are estimated by using the ordinary least square or linear least
264 square method. This method tries to minimise the sum of the squares of the error terms.

265 MLR model can be evaluated using statistical tools. The adjusted coefficients of determination (R^2_{adj}), is a statistical
 266 index that provides information about the goodness of fit of a model. It represents the proportion of the variation in the
 267 dependent variable that is attributable to the explanatory variables:

$$268 \quad R^2_{adj} = 1 - \left[\frac{(1-R^2)(n-1)}{(n-p-1)} \right] \quad (8)$$

269 where R^2 is the coefficient of determination, n is the number of observations, and p is the number of variables included
 270 in the model.

271 The coefficient t student is used to test the null hypothesis, i.e. when the values of the estimated coefficients of MLR
 272 model are not significant:

$$273 \quad t = \beta_{1,\dots,p} / SE_{\beta_{1,\dots,p}} \quad (9)$$

274 where $\beta_{1,\dots,p}$ are the estimated coefficients of MLR model and $SE_{\beta_{1,\dots,p}}$ is the standard error of each the estimated
 275 coefficient. Generally, if, $t \leq |2|$, $\beta_{1,\dots,p}$ is less significant.

276 A method for testing the significance of the MLR model is the Fisher-Snedecor test (F). It is conducted on the entire
 277 model and is based on the decomposition of deviance:

$$278 \quad F = MS / MS_R \quad (10)$$

279 where MS is the Means Square of the model and MSr is the Residual Mean Square. If the value of F does not exceeds
 280 the critical value (default value for a given probability), the correlation between the variables is not linear. Therefore,
 281 there could be a different correlation.

282 Durbin-Watson ($D-W$) is a statistic test used to detect the presence of autocorrelation in the residuals (estimation
 283 errors) of a MLR model.

$$284 \quad D-W = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \quad (11)$$

285 where $e_i = y_i - \hat{y}_i$ and y_i and \hat{y}_i are respectively the observed values and the expected value of the response variable for
 286 individual i . The value of D-W always lies between 0 and 4. As a rule of thumb, $D-W = 2$ indicates no autocorrelation,
 287 $D-W < 2$ indicates negative autocorrelation, and $D-W > 2$ indicates positive autocorrelation.

288 **2.4 Classification and Regression Tree method**

289 The CART [34, 35] is a binary decision tree that is constructed by splitting a parent node into two child nodes
 290 repeatedly, beginning with the root node that contains the whole learning sample. The CART can easily handle both
 291 numerical and categorical variables.

292 A decision tree generation consists of a two-step process: learning and classification. In the first step, the dataset is
 293 divided into a training set and a testing set. The creation of these two subsets is the most delicate part of the technique.
 294 It is important that the training set and the testing set come from the same population and that they are disjointed. In the
 295 classification process, the results obtained from the training set are the input to test the decision tree. The accuracy of
 296 the model is measured by comparing the estimated values of each "leaf node" with the real values contained in the test
 297 sample. If the estimation is acceptable, the decision tree can be applied to new datasets for classification and estimation.
 298 Initially, all records in the training data are grouped together into a single unit. At each iteration, the algorithm chooses

299 a predictor attribute that can “best” separate the target class values. Measures of impurity are used to estimate the ability
 300 of a predictor to separate the target class values. The CART consists of three parts:

301 1. Construction of maximum tree: the classification tree is built in accordance with the splitting rule. Each time
 302 data must be divided into two parts with the highest homogeneity. The Least Square Deviation measures the
 303 impurity of a node t and is defined as:

$$304 \quad i(T) = \frac{1}{N_w(t)} \cdot \sum_{i=t} w_i \cdot f_i \cdot (y_i - \bar{y}(t))^2 \quad (12)$$

305 where N_w is the weighted number of cases in node t , w_i is the weight of the variable in the case, f_i is the
 306 frequency value of the variable, y_i is the value of the response variable, and $\bar{y}(t)$ is the weighted average value
 307 of the variable at node t . The best split s^* of a generic node t is what determines the greater decrease of the $i(t)$.
 308 For each split s of node t into t_r and t_l the following algorithm is valid:

$$309 \quad \Delta i(s^*, t) = i(T) - (i(t_l) + i(t_r)) \quad (13)$$

310 where $\Delta i(s^*, t)$ is the decrease of impurities in a generic node, $i(T)$ is the Least Square Deviation of the given
 311 node, $i(t_{l,r})$ are the Least Square Deviation of the two nodes split.

312 2. Choice of the right tree size. Two rules for the stop can be used in practice: optimisation by number of points in
 313 each node (minimum number of cases in the parent’s node and the child’s node) and the error E associated with
 314 the tree. The E parameter allows the tree to be properly built:

$$315 \quad E = \sum_{i=1}^n \left(n_i / n_i^* \cdot \varepsilon_i \right) \quad (14)$$

316 where n_i is the total number of records in the training set which terminate in the leaf, and n_i^* is the number of
 317 records classify bin the leaf i .

318 The optimal condition is obtained by setting the error $E=1$. In fact, in this case the tree correctly classifies all of
 319 the records in the training set. The optimisation of the tree size is important, because the maximum trees may
 320 turn out to be very complex and may consist of hundreds of levels.

321 3. Classification of new data: each of the new observations will be set to one of the terminal nodes of the tree by
 322 means of a set of questions. A new observation is assigned with the dominating class/response value of the
 323 terminal node, where this observation belongs.

324 3. Results: development of models

325 3.1 Outliers detection analysis

326 In order to find potential outliers, a pre-processing analysis was carried out prior to creating the estimation model. All of
 327 the variables should show a sufficient range of variability and have skewness and kurtosis values of less than $|1.00|$.
 328 Indeed, including variables whose distribution is too different from the normal value into the MLR model can lead to
 329 the violation of the assumptions of linearity and homoskedasticity of the residual anomalies. The variable *floor heated*
 330 *area* was excluded, because its values were missing for 18 schools.

331 The pre-processing analysis identified 14 potential outliers. After conducting an accurate frequency distribution analysis
 332 of the sample, it was observed that the detected outliers belonged to the tails of distribution for each variable. In
 333 particular, it was verified that these outliers influence the mean and standard deviation for each variable, causing a non-
 334 normal distribution for all of them. Even if the detected outliers can be considered reliable from an energy measurement
 335 point of view, it was verified that they decrease the performance of the estimation models. A detailed analysis on these

336 buildings showed that they are characterised by high thermal transmittance values, low system energy efficiencies, low
 337 number of pupils, very low or very high volume. This is the reason why the outliers belong to the tails of distribution
 338 for each variable and therefore determine anomalous values of heating energy consumption. Table 2 shows the indexes
 339 evaluated for the sample without outliers.

340 3.2 Multiple Linear Regression model

341 In order to develop the MLR model, all of the reported anomalies were deleted by excluding the identified outliers from
 342 the database. The assumptions relating to the specification of the model (do not omit relevant predictors and do not
 343 include irrelevant predictors) were verified by evaluating the bivariate correlations between the independent variables
 344 and the dependent variable constituted by heating energy consumption (Figure 6 and Table 3).

345 The variables *building height* and *number of floors* have a correlation coefficient of less than 0.20 (Figure 6). For this
 346 reason, they were not included in the model.

347 The values of the parameter Variance Inflation Factors (VIF) are reported in Table 3. VIF allows detecting the presence
 348 of multicollinearity between the explanatory variables. In general, a multicollinearity occurs if the value of VIF exceeds
 349 10.

350 The results analysis found a strong correlation between:

- 351 - number of pupils and number of classrooms;
- 352 - aspect ratio and gross heated volume;
- 353 - aspect ratio and heat transfer surface.

354 Given these findings, the variables included in the MLR model were reduced to: *real heating degree days*, *gross heated*
 355 *volume*, *heat transfer surface*, *thermal transmittance of walls and windows*, *boiler size*, *number of pupils*, *annual*
 356 *operating time*, and *average seasonal system efficiency*, as summarised in Table 4.

357 The sample was randomly split into the training dataset (39 records were selected from the database, i.e. 70 % of the
 358 sample) and testing dataset (the remaining 27 records, i.e. 30 % of the sample). The estimation model was therefore
 359 developed on the basis of a training sample. The training set does not include the outliers previously identified. Each
 360 variable was standardised by the z-score method (Eq.1) to compare variables between them by assuming the same
 361 distribution ($\mu = 0$; $\sigma = 1$). The most accurate estimation for heating energy consumptions (measured in kWh) is
 362 calculated by means of the following equation:

$$363 \quad EUI_{st} = 662765 + \beta_1 \cdot X_1^* + \beta_2 \cdot X_2^* + \beta_3 \cdot X_3^* - \beta_4 \cdot X_4^* - \beta_5 \cdot X_5^* - \beta_6 \cdot X_6^* + \beta_7 \cdot X_7^* + \beta_8 \cdot X_8^* - \beta_9 \cdot X_9^* \quad (15)$$

364 where the variables of the model are shown in Table 5, including detailed information about the estimated coefficients
 365 (β) of MLR model, partial standardised regression coefficients (b), and the t-values.

366 All of the examined variables within this study can theoretically impact heating demand. The *gross heated volume*,
 367 *boiler size*, and *thermal transmittance of windows* exhibit the greatest impact in the model, with partial standardised
 368 regression coefficients of 0.86, 0.64, and 0.61, respectively. The t-test identifies the inference on individual coefficients
 369 β . In particular, it verifies whether every single variable X^* influences the response variable. The variable *annual*
 370 *operating time* and *average seasonal system efficiency* are the only two variables with a t-value of less than |2|. For this
 371 reason, both of their estimated coefficients of MLR model are less significant. The variance showed by the model

372 compared to the total variance of the sample is 86 % (R^2_{adj}), therefore, 86 % of the heating energy consumption variance
 373 can be explained by the nine variables used in the model. Moreover, there is an absence of auto-correlations among
 374 residuals ($D-W = 2.05 \approx 2$) and the value of the Fisher-Snedecor test ($F = 27$) is greater than the critical value ($F_{crit} = 4$).
 375 As such, the MLR model can be considered robust. To assess the quality of the estimation model, Figure 7 shows the
 376 distribution plot between the estimated EUI_{st} and the monitored EUI_{st} using the testing dataset.

377 The best fit is affected by an error of 1 %, while the bad fit of the model is affected by an error of 40 %. The average
 378 error is equal to 15 % and for testing dataset the value of R^2 is 86 %. The model tends to underrate the energy
 379 consumption; in fact, 16 cases of 27 show an estimated heating energy consumption that is lower than the actual value.
 380 The validation test demonstrates that the model has an adequate estimation ability.

381 3.3 Classification and Regression Tree

382 In order to develop the CART, the training dataset and the test dataset used are the same as the ones used in the
 383 regression model. The CART algorithm selected five parameters from the database to model input variables (Table 6).

384 The decision tree was constructed to estimate the heating energy consumptions. The rules set for the arrest of the tree
 385 are as follows:

- 386 – minimum number of cases (*parents node*): 2
- 387 – minimum number of cases (*children node*): 2
- 388 – $E = 1$.

389 The tree includes a total of eight leaf nodes that represent the final classes. The estimation of the heating energy
 390 consumptions of each leaf node corresponds with the average of cases included in it. The algorithm can be translated
 391 into a set of decision rules that take the following form: *if* antecedent conditions, *then* consequent conditions. Table 7
 392 presents the results of the CART in terms of the decision rules for the training dataset, starting from the root node and
 393 following all the way to each leaf node.

394 The decision rules can be used to estimate the EUI_{st} target level of a new school building having similar features. For
 395 example, looking at the first rules (Table 7), the EUI_{st} level can be estimated as follows:

396 Step 1: The root node is the starting point for the estimation. Table 7 shows that the value of the VOL variable should
 397 be examined first. If the VOL is higher or equal to 33195 m^3 , then it is possible to go to the next step.

398 Step 2: examine the value of the SUR variable; if SUR is lower or equal to 12818 m^2 , the EUI_{st} level of the school
 399 building is 968 MWh.

400 The CART carries out a sensitivity analysis before creating the decision rules in order to select the variables more
 401 correlated with the heating energy consumptions. In fact, the variables selected by the algorithm are characterised by
 402 high correlation coefficients (see Figure 6). The other factors (heating degrees days, thermal transmittance of walls,
 403 number of pupils and classrooms, average seasonal system efficiency, and annual operating time) do not appear in the
 404 decision tree, because they were excluded during the pruning process.

405 As previously mentioned, the accuracy of the decision tree must be evaluated before it is applied to a new dataset. Since
 406 the estimated values correspond to the mean value of the data included in the node, the estimation will always be
 407 affected by an error. For this reason, it is appropriate to associate a confidence interval for each estimated value. The
 408 confidence intervals with a 95 % probability of containing the true parameters were calculated and the results are shown
 409 in Table 8.

410 The decision tree was applied to the testing dataset and the results are reported in Figure 8. The best fit is affected by an
 411 error of 2 %, while the model's bad fit is affected by an error of 33 %. The average error is equal to 13 %, and the value
 412 of R^2 for the testing dataset is 86 %.

413 3.4 Models comparison

414 Two estimation models were evaluated based on their ability to estimate heating energy consumption in school
 415 buildings. Although they share the same goal, these models are based on different methodologies. To understand the
 416 possibilities and limitations of both, a residuals analysis is needed. The residuals are often used as an indicator to
 417 validate a model. The basic hypothesis is that the residuals, i.e. the errors, are randomly distributed. Moreover, they
 418 should not be correlated with the dependent or independent variables and the average value of the residuals should be
 419 equal to zero. The last hypothesis was verified for both of the models. The value is less uncertain for the residuals of
 420 CART than for the MLR model. No significant correlation was identified between the variables in both of the models.
 421 Comparing the goodness of fit of the two models, in percentage terms, MLR model commits an average error of 15 %,
 422 while the CART commits an average error equal to 13 %. In particular, the variance explained by the two methods is
 423 equal to 86 %. The residuals analysis and the coefficient of determination are not enough to evaluate the performance of
 424 the estimation models. Three other criteria have been used to test the performance of both models, the Mean Absolute
 425 Error (MAE) the Root Mean Square Error (RMSE) and the Mean Absolute Percentage error (MAPE). These parameters
 426 are shown in the following equations:

$$427 \quad MAE(N) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (16)$$

$$428 \quad RMSE(N) = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}} \quad (17)$$

$$429 \quad MAPE(N) = \frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right| \cdot 100 \quad (18)$$

430 where y_i is the monitored heating energy consumptions and \hat{y}_i is the estimated heating energy consumptions and N is
 431 the sample size.

432 Table 9 gathers the results of the error analysis. The three indexes for the CART are always lower than the values
 433 estimated for the MLR model. The MAPE values are very similar for both models, but the MAE and the RMSE,
 434 considering the same testing dataset, are lower for the CART. A low RMSE value means that the error is characterised
 435 by a low dispersion. This is more clear in Figure 9, where a comparison between measured and estimated EUI_{st} is
 436 reported with a box-plots representation.

437 The range of the measured EUI_{st} is comparable with the estimated outputs of both models. In fact, the median values for
 438 the three bars is equal. The lower and upper quartile values of the measured data and the data estimated with CART are
 439 quite close. This confirms that the MLR model tends to underrate the heating energy consumption (minimum, median,
 440 and lower quartile are closer compared to the monitored data), while the CART's output is unbiased. The results show a
 441 strong relationship between the dependent variables and the heating energy consumptions. In particular, the most
 442 influential variables in the MLR model (gross heated volume, heat transfer surface, boiler size, and thermal
 443 transmittance of windows) are exactly the same factors that CART selected to create the decision rules. It can be

444 concluded that both of the estimation models were correctly developed. CART's performance is slightly better than that
445 of MLR model, as demonstrated by the results presented earlier (residual analysis, R^2 , MAE, RMSE, and MAPE).

446 **4. Discussion**

447 The analysis of CRESME [2] shows that the Italian school building stock could achieve energy savings of about 48.3 %
448 and shift from a current energy consumption rate of 9.6 TWh/yr to a target value of 5.0 TWh/yr. Similar results may be
449 obtained for public directional buildings and dwellings. However, a more detailed analysis may be conducted at the
450 local level, enabling a better definition of actual planning actions and economic assessments.

451 The actions that a local authority may adopt to challenge energy savings in the construction sector include:

- 452 - defining a public building portfolio (building stocks) and reference performance benchmarks;
- 453 - setting simple thresholds for energy performance, using existing energy data;
- 454 - setting a priority action list for the energy management and renovation of the building portfolio;
- 455 - adopting economic policies to promote the most relevant actions.

456 These actions may be part of an effective energy plan, but they require a set of technical steps: all of the considered
457 asset's fundamental data must be collected, gathered, and assembled in an appropriate database, eventual outliers must
458 be processed, and an analysis then ultimately produces effective decision rules. These include both the planning of
459 ordinary management activities and strategic planning targeting energy efficiency improvements.

460 The choice of the most adequate and accurate estimation model to perform the required analysis, knowing its
461 possibilities and limitations, is crucial in order to correctly inform the following local authority actions.

462 The estimation accuracy of the two models analysed in the present paper showed to be influenced by the nature of the
463 dataset, in terms of its density and how uniform the frequency distribution is. The CART is based on binary splitting
464 criteria of the response variable as a function of the influencing variables. It performs well when the leaf nodes are
465 characterised by values that are close to the mean (low confidence intervals). In this case also numerical variables can
466 be used as target attribute. However, generally the CART algorithm is used to classify categorical attribute. On the other
467 hand, a non-uniform dataset could make the MLR model incapable of estimating unbiased regression coefficients.

468 The MLR model requires knowing the exact values of all input variables. This is a weakness, in fact, as the precise
469 value of some of the variables, such as the thermal transmittance of walls/windows or heat transfer surface, is not easily
470 obtained or readily available for existing buildings. Moreover, the model requires the input parameters to be
471 standardised. For this reason, it is not easy for inexperienced users to interpret and use. On the other hand, the MLR
472 model can be used to create benchmarks [36]. A benchmark value may be used as a target to be reached or exceeded
473 and may prove quite useful to guide designers towards the optimal technical and economical solution. This is not
474 possible with the CART.

475 Despite being a particular data mining technique, the CART's output consists of a set of decision rules that even non-
476 experts can easily understand and use. Useful information can be obtained from this model, for example, it helps to
477 understand a building's energy consumptions pattern and how to optimise a building's design. The algorithm
478 automatically selects the different parameters as predictors. These are used to split the nodes of the decision tree, and
479 their proximity to the root node indicates the strength of the influence and the number of records impacted. By
480 examining the decision rules (see table 7), one can identify what primary factors account for the energy demand profiles
481 of the schools. Among the considered factors, the root node, i.e. gross heated volume, indicates that the size of the

482 schools is the most important element in determining energy demand. The heating energy consumption of large school
483 buildings (Rules 1 – 2) is only influenced by the gross heated volume and heat transfer surface. Instead, in medium size
484 school buildings (Rules 6 – 7 – 8), the heating energy consumption is influenced by four significant factors (gross
485 heated volume, heat transfer surface, boiler size, and thermal transmittance of windows). Finally, the heating energy
486 consumption of small size school buildings (Rules 3 – 4 – 5) is a function of two geometric factors (gross heated
487 volume and heat transfer surface) and one construction feature (thermal transmittance of windows).

488 The accuracy of MLR and CART models results quite good also compared with values obtained by other researches.
489 For example in [16], eight different regression models (Linear – Logarithmic – Quadratic – Cubic – Inverse – Linear
490 and Inverse – Power – S – Exponential) were developed to estimate the energy consumptions of 105 schools located in
491 Germany. The linear and inverse regression model showed the best fit with a MAPE of 17 %. Thewes et al. [17]
492 developed a regression model able to explain 53 % ($R^2 = 53\%$) of electrical and heating energy consumption variance
493 for 68 school buildings situated in Luxembourg. In [15], a multiple regression model was developed to estimate daily
494 electricity consumption of an administration building located at the Southwark campus of London South Bank
495 University in London. The final model has an adjusted R^2 value of 88 %.

496 No example of the CART model used for the estimation of heating energy consumption in schools is available.
497 Therefore, only results obtained for other building types may be reported. Yu et al. [18], used a decision tree and
498 decision rules to classify the EUI level of a new residential building in Japan. The C4.5 algorithm was used with a
499 percentage error between 0.2 and 141.9 % (average error equal to 25 %), on the basis of 55 records for the training
500 dataset and 12 records for the test dataset.

501 **5. Conclusion**

502 The present work studied two estimation models, based on different modelling methodologies, and applied them to
503 estimate the heating consumption of school buildings in the north of Italy. The methods compared are: a multiple linear
504 regression model and a classification and regression tree. While MLR model have been successfully applied in former
505 works, data mining techniques, such as the decision tree, are a newly emerging analysis tool. The application of the
506 decision tree to school buildings was demonstrably reliable in terms of the heating energy consumption estimation. The
507 variance explained by both models is 86 %, but the decision tree shows lower errors, evaluated by means of the MAE,
508 RMSE, and MAPE. Moreover, the gross heated volume, heat transfer surface, boiler size, and thermal transmittance of
509 windows, were the parameters identified as primarily influencing the heating energy consumption among the considered
510 school building stock.

511 The two methods are complimentary, not antagonistic, and show different strengths and weaknesses, as discussed in this
512 paper. The greatest advantage of the CART is that the output consists of a set of practical decision rules that decision
513 makers can quickly use. It also provides useful information on the influencing variables for each leaf node representing
514 a sub-dataset, i.e. a homogenous class of school buildings.

515 The MLR model output consists of an equation including all of the major variables affecting heating energy
516 consumption. Moreover, the partial standardised regression coefficients provide information on the most influencing
517 input variables, making it possible to carry out a sensitivity analysis. Finally, the MLR model can be used to perform
518 benchmark analyses.

519 Since the variability of the analysed sample is large enough to represent all school buildings in the north of Italy, the
520 developed models can be used to estimate the heating energy requirements of new structures whose characteristics are

521 within the ranges (for each variable) reported for the training dataset of the models. Similar models may moreover be
 522 developed, on the basis of other database, to estimate the heating energy requirements of different building types, since
 523 the model showed to be reliable and robust.

524 Future research should concern the possibility to couple the two models, in order to increase further the estimation
 525 capability. When the data of each final class of the CART are characterised by a large confidence interval, the
 526 performance of the model decreases. For each of these nodes is therefore possible to develop a MLR model, since each
 527 node is made by a sub-dataset of buildings with homogenous features. This combination may exploit the best
 528 characteristic of each of the two models; however, it requires new and larger training and testing database.

529 REFERENCES

- 530 [1]L. Perez-Lombard, J. Ortiz, C. Pout, A review on buildings energy consumption information, *Energy and Building* 40 (2008) 394-
 531 398.
- 532 [2] CRESME, RIUSO03: Ristrutturazione edilizia riqualificazione energetica rigenerazione urbana – estratto della ricerca CRESME,
 533 2014 (Italian text).
- 534 [3]EU Directive 2002/91/EC of the European Parliament and of the Council of 16 December 2002 on the energy performance of
 535 buildings.
- 536 [4]EU Directive 2010/31/EU of the European Parliament and of the Council of 19 May 2010 on the energy performance of buildings.
- 537 [5]EU Directive 2012/27/EU of the European Parliament and of the Council of 25 October 2012 on energy efficiency, amending
 538 Directives 2009/125/EC and 2010/30/EU and repealing Directives 2004/8/EC and 2006/32/EC.
- 539 [6]D. Fiaschi, R. Bandinelli, S. Conti, A case study for energy issues of public buildings and utilities in a small municipality:
 540 Investigation of possible improvements and integration with renewables, *Appl. Energy*. 97 (2012) 101–114.
- 541 [7]D. Radulovic, S. Skok, V. Kirincic, Energy efficiency public lighting management in the cities, *Energy*. 36 (2011) 1908–1915.
- 542 [8]US DOE-Energy Smart Schools, US Department of Energy, Energy Efficiency and Renewable Energy, Building Technologies
 543 Program. Available from: <<http://www.eere.energy.gov/buildings/energysmartschools/resources.html>>.
- 544 [9]L. De Santoli, F. Fraticelli, F. Fornari, C. Calice, Energy performance assessment and a retrofit strategies in public school
 545 buildings in Rome, *Energy Build.* 68 (2014) 196–202.
- 546 [10]A. Dimoudi, P. Kostarela, Energy monitoring and conservation potential in school buildings in the C' climatic zone of Greece,
 547 *Renew. Energy*. 34 (2009) 289–296.
- 548 [11]T.W. Kim, K.G. Lee, W.H. Hong, Energy consumption characteristics of the elementary schools in South Korea, *Energy Build.*
 549 54 (2012) 480–489.
- 550 [12]S.P. Corgnati, V. Corrado, M. Filippi, A method for heating consumption assessment in existing buildings: A field survey
 551 concerning 120 Italian schools, *Energy Build.* 40 (2008) 801–809.
- 552 [13]S.P. Corgnati, F. Ariaudio, L. Rollino, Definizione di un indice semplificato per la previsione dei consumi per il riscaldamento di
 553 un patrimonio edilizio esistente a destinazione d'uso prevalentemente scolastica, III Congresso Nazionale AIGE, Parma 4-5 June
 554 2009.
- 555 [14]F. Ariaudio, S.P. Corgnati, M. Filippi, Heating consumption assessment and forecast of existing buildings: investigation on
 556 Italian school Buildings, *Proceedings of the 5th IBPC, KYOTO* (2012) May 28-31.
- 557 [15]K.P. Amber, M.W. Aslam, S.K. Hussain, Electricity consumption forecasting models for administration buildings of the UK
 558 higher education sector, *Energy Build.* 90 (2015) 127–136.
- 559 [16]E. Beusker, C. Stoy, S.N. Pollalis, Estimation model and benchmarks for heating energy consumption of schools and sport
 560 facilities in Germany, *Build. Environ.* 49 (2012) 324–335.
- 561 [17]A. Thewes, S. Maas, F. Scholzen, D. Waldmann, A. Zürbes, Field study on the energy consumption of school buildings in
 562 Luxembourg, *Energy Build.* 68 (2014) 460–470.
- 563 [18]Z. Yu, F. Haghghat, B.C.M. Fung, H. Yoshino, A decision tree method for building energy demand modeling, *Energy Build.* 42
 564 (2010) 1637–1646.

- 565 [19]G.K.F. Tso, K.K.W. Yau, Predicting electricity energy consumption: A comparison of regression analysis, decision tree and
566 neural networks, *Energy*. 32 (2007) 1761–1768.
- 567 [20]A.Z. Al-Garni, S.M. Zubair, J.S. Nizami, A regression model for electric-energy consumption forecasting in Eastern Saudi
568 Arabia. *Energy* 19 (1999) 1043-1049.
- 569 [21]A. Aranda, G. Ferreira, M.D. Mainar-Toledo, S. Scarpellini, E. Llera Sastresa, Multiple regression models to predict the annual
570 energy consumption in the Spanish banking sector, *Energy Build.* 49 (2012) 380–387.
- 571 [22]I. Korolija, Y. Zhang, L. Marjanovic-Halburd, V.I. Hanby, Regression models for predicting UK office building energy
572 consumption from heating and cooling demands, *Energy Build.* 59 (2013) 214–227.
- 573 [23]R.Z. Freire, G.H.C. Oliveira, N. Mendes, Development of regression equations for predicting energy and hygrothermal
574 performance of buildings, *Energy Build.* 40 (2008) 810–820.
- 575 [24]J. Zhao, B. Lasternas, K.P. Lam, R. Yun, V. Loftness, Occupant behavior and schedule modeling for building energy simulation
576 through office appliance power consumption data mining, *Energy Build.* 82 (2014) 341–355.
- 577 [25]R. Mikučionienė, V. Martinaitis, E. Keras, Evaluation of energy efficiency measures sustainability by decision tree method,
578 *Energy Build.* 76 (2014) 64–71.
- 579 [26]L. Dias Pereira, D. Raimondo, S.P. Corgnati, M. Gameiro da Silva, Energy consumption in schools – A review paper, *Renew.*
580 *Sustain. Energy Rev.* 40 (2014) 911–922.
- 581 [27]S.P. Corgnati, E. Fabrizio, F. Ariaudo, L. Rollino: Edifici tipo, indici di benchmark di consumo per tipologie di edificio, ad uso
582 scolastico (medie superiori e istituti tecnici) applicabilità di tecnologie innovative nei diversi climi italiani, Report RSE/2010/190.
- 583 [28]N. Gaitani, C. Lehmann, M. Santamouris, G. Mihalakakou, P. Patargias, Using principal component and cluster analysis in the
584 heating evaluation of the school building sector, *Applied Energy*. 87 (2010) 2079–2086.
- 585 [29]R.E. Walpole, R.H. Myers, S.L. Myers, K. Ye., *Probability & Statistics for Engineers & Scientists*, ninth ed., Academic Press,
586 New York, 2011.
- 587 [30]M. Heidarinejad, M. Dahlhausen, S. McMahon, C. Pyke, J. Srebric, Cluster analysis of simulated energy use for LEED certified
588 U . S. office buildings, *Energy Build.* 85 (2014) 86–97.
- 589 [31]F. Xiao, C. Fan, Data mining in building automation system for improving building operational performance, *Energy Build.* 75
590 (2014) 109–118.
- 591 [32]J. Wakefield, *Bayesian and Frequentist Regression Model*, Springer-Verlag, New York, 2012.
- 592 [33]S.J. Sheather, *A Modern Approach to Regression with Regression Model*, Springer-Verlag, New York, 2009.
- 593 [34]C. Charu Aggarwal, *Data Classification Algorithm and Application*, Chapman & Hall/CRC, New York, 2014.
- 594 [35]L. Breiman, J. Friedman, C.J. Stone, R.A. Olshen, *Classification and Regression Trees*, Chapman & Hall/CRC, New York,
595 1984.
- 596 [36]W. Chung, Y. V. Hui, Y.M. Lam, Benchmarking the energy efficiency of commercial buildings, *Applied Energy*. 83 (2006) 1–
597 14.

Variables	Minimum	Maximum	Mean	Standard Deviation
Real Heating Degree Days [DD]	2537	3197	2701	161
Gross Heated Volume [m ³]	2900	86830	31464	20010
Heat Transfer Surface [m ²]	1905	25946	9866	5404
Aspect Ratio [m ⁻¹]	0.20	0.80	0.35	0.10
Floor Heated Area [m ²]	1578	18254	6791	4073
Building Height [m]	5	28	14.50	4.50
Numbers of Floors [n°]	2	7	3	1.5
Thermal Trasmittance of Walls [W/m ² K]	0.40	2.39	1.05	0.50
Thermal Trasmittance of Windows [W/m ² K]	2.90	6.50	4.30	0.90
Boiler Size (Heat input) [kW]	106	8000	1755	1345
Number of Classrooms [n°]	6	65	27	14
Number of Pupils [n°]	75	1340	583	320
Annual Operating Time [h]	889	1888	1444	2450
Average Seasonal System Efficiency [%]	0.45	0.86	0.64	0.07

Tab 1 Statistical description of the variables influencing the heating energy consumption

Parameters	Limit Values	Calculated Values
z-score	< 3	No Values Higher
Mahalanobis Distance	27.87	No Values Higher
Index of Mardia	99	94.85
Cook Distance	1	0.29
Leverage Values	0.60	No Values Higher

Tab. 2 Indexes for the outliers detection analysis

Accepted Manuscript

Variables	VIF
Real Heating Degree Days [DD]	1.53
Gross Heated Volume [m ³]	24.28
Heat Transfer Surface [m ²]	19.25
Aspect Ratio [m ⁻¹]	4.96
Building Height [m]	2.55
Numbers of Floors [n ^o]	2.46
Thermal Trasmittance of Walls [W/m ² K]	1.54
Thermal Trasmittance of Windows [W/m ² K]	1.13
Boiler Size (Heat input) [kW]	3.10
Number of Classrooms [n ^o]	35.90
Number of Pupils [n ^o]	36.75
Annual Operating Time [h]	2.02
Average Seasonal System Efficiency [%]	1.347

Tab. 3 Variance Inflation Factors (VIF)

Accepted Manuscript

Variables	Minimum	Maximum	Mean	Standard Deviation
Real Heating Degree Days [DD]	2537	3113	2696	159
Gross Heated Volume [m ³]	5065	78532	28745	17497
Heat Transfer Surface [m ²]	1905	24206	9354	5122
Thermal Trasmittance of Walls [W/m ² K]	0.40	2.39	1.05	0.40
Thermal Trasmittance of Windows [W/m ² K]	2.90	6.50	4.25	0.95
Boiler Size (Heat input) [kW]	141	3807	1424	845
Number of Pupils [n°]	115	1194	530	266
Annual Operating Time [n°]	889	1848	1426	250
Average Seasonal System Efficiency [%]	0.45	0.77	0.64	0.06

Tab 4 Statistical description of the variables included into MLR model

	Variables	β	b	t-value
	Intercepts	662757	-	3.82
X ₁ [*]	Z: Real Heating Degree Days	16954	0.36	2.11
X ₂ [*]	Z: Gross Heated Volume	203734	0.86	6.29
X ₃ [*]	Z: Heat Transfer Surface	22609	-0.44	2.49
X ₄ [*]	Z: Thermal Trasmittance of Walls	-12285	-0.39	-2.44
X ₅ [*]	Z: Thermal Trasmittance of Windows	-24142	-0.61	-2.89
X ₆ [*]	Z: Boiler Size (Heat input)	103167	-0.64	3.40
X ₇ [*]	Z: Number of Pupils	55185	0.51	2.39
X ₈ [*]	Z: Annual Operating Time	9940	0.03	0.35
X ₉ [*]	Z: Average Seasonal System Efficiency	-1131	-0.04	-0.55

Tab. 5 Estimated coefficients (β) partial standardized regression coefficients (b) and t-values

Accepted Manuscript

N°	VARIABLE	NAME	TYPE	Unit of measure
1	Standard Energy Use Intensity	EUI _{ST}	Numerical	MWh
2	Heated Gross Volume	VOL	Numerical	m ³
3	Heat Transfer Surface	SUR	Numerical	m ²
4	Boiler Size (Heat input)	POW	Numerical	kW
5	Thermal Transmittance of Windows	U _{windows}	Numerical	W/m ² K

Tab 6 Variables selected in the CART

Accepted Manuscript

N° RULES	DECISION RULES
1	If $VOL \geq 33195 \text{ m}^3$ and SUR is $< 12818 \text{ m}^2$ then EUI_{st} is 968 MWh
2	If $VOL < 33195 \text{ m}^3$ and SUR is $\geq 12818 \text{ m}^2$ then EUI_{st} is 1183 MWh
3	If $VOL < 33195 \text{ m}^3$ and $SUR < 2460 \text{ m}^2$ then EUI_{st} is 140 MWh
4	If $VOL < 33195 \text{ m}^3$ and $SUR < 6203 \text{ m}^2$ and $SUR \geq 2460 \text{ m}^2$ and $U_{windows}$ is $< 4.65 \text{ W/m}^2\text{K}$ then EUI_{st} is 303 MWh
5	If $VOL < 33195 \text{ m}^3$ and $SUR < 6203 \text{ m}^2$ and $SUR \geq 2460 \text{ m}^2$ and $U_{windows}$ is $\geq 4.65 \text{ W/m}^2\text{K}$ then EUI_{st} is 421 MWh
6	If $VOL < 33195 \text{ m}^3$ and SUR is $\geq 6203 \text{ m}^2$ and $U_{windows}$ is $< 4.54 \text{ W/m}^2\text{K}$ then EUI_{st} is 521 MWh
7	If $VOL < 33195 \text{ m}^3$ and SUR is $\geq 6203 \text{ m}^2$ and $U_{windows}$ is $\geq 4.54 \text{ W/m}^2\text{K}$ and POW is $< 1336 \text{ kW}$ then EUI_{st} is 708 MWh
8	If $VOL < 33195 \text{ m}^3$ and SUR is $\geq 6203 \text{ m}^2$ and $U_{windows}$ is $\geq 4.54 \text{ W/m}^2\text{K}$ and POW is $\geq 1336 \text{ kW}$ then EUI_{st} is 816 MWh

Tab 7 Decision rules

Accepted Manuscript

RULES	ESTIMATED EUI_{st} [MWh]	UPPER CONFIDENCE LEVEL [MWh]	LOWER CONFIDENCE LEVEL [MWh]
1	968	1004	931
2	1183	1304	1061
3	140	214	56
4	303	328	279
5	421	472	369
6	521	604	439
7	708	774	641
8	816	832	801

Tab 8 Confidence interval of estimated EUI_{st} (CART)

Accepted Manuscript

INDEX	MLR MODEL	CART
MAE (MWh)	108	102
RMSE (MWh)	145	142
MAPE (%)	15	14

Tab 9 Error comparison for MLR model and CART

Accepted Manuscript

Fig.1 Framework of the research

Fig.2 Sample description: gross heated volume and external walls surface

Fig.3 Sample description: thermal transmittance of walls and windows

Fig.4 Sample description: boiler size (heat input) and average seasonal system efficiency

Fig.5 Sample description: heating energy consumption

Fig.6 Correlation coefficients

Fig. 7 Distribution plot between the monitored and the estimated EUI_{st} (MLR Model testing dataset)

Fig. 8. Distribution plot between the monitored and the estimated EUI_{st} (CART - testing dataset)

Fig. 9. Box-plots of monitored and estimated heating energy consumption

Figure 1
[Click here to download high resolution image](#)

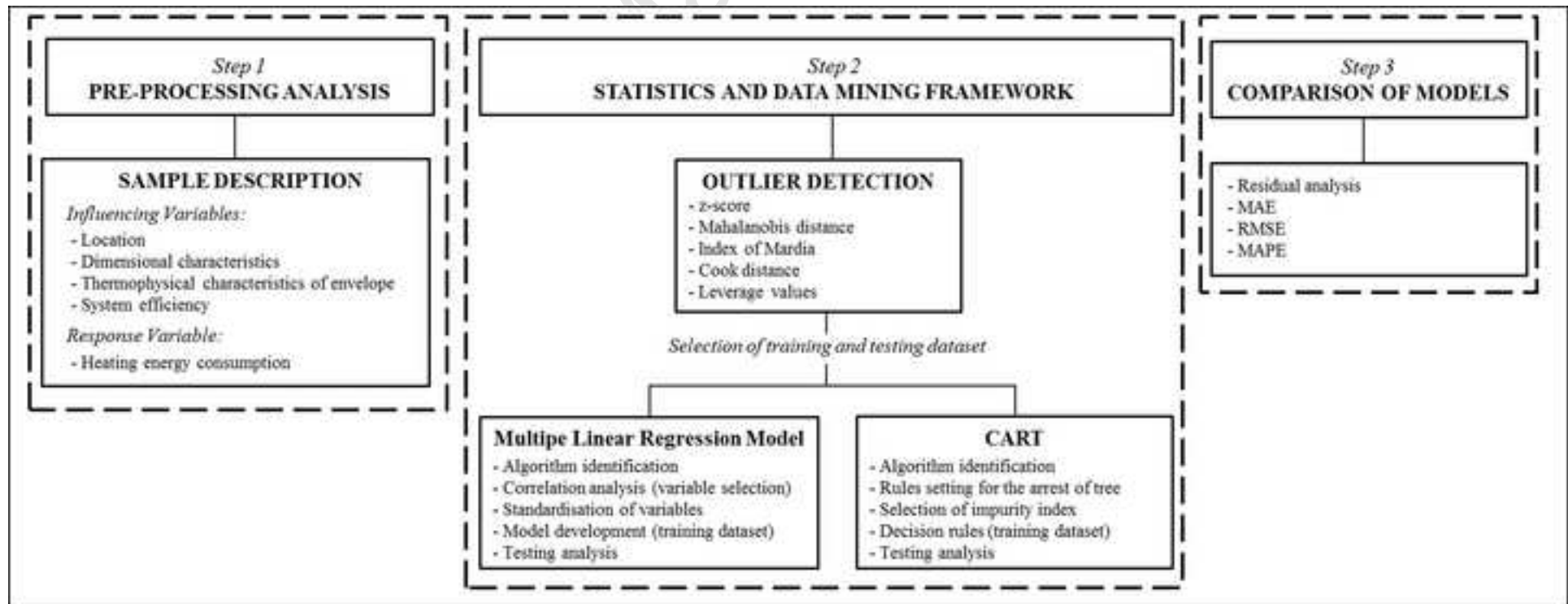


Figure 2
[Click here to download high resolution image](#)

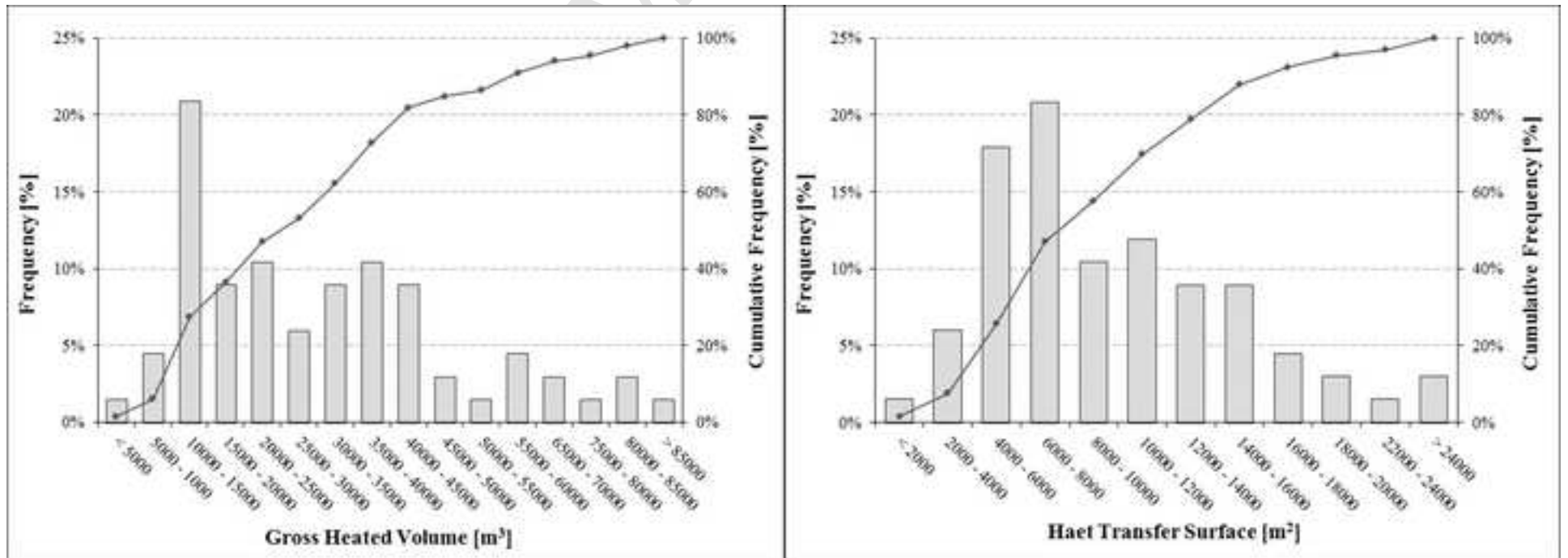


Figure 3
[Click here to download high resolution image](#)

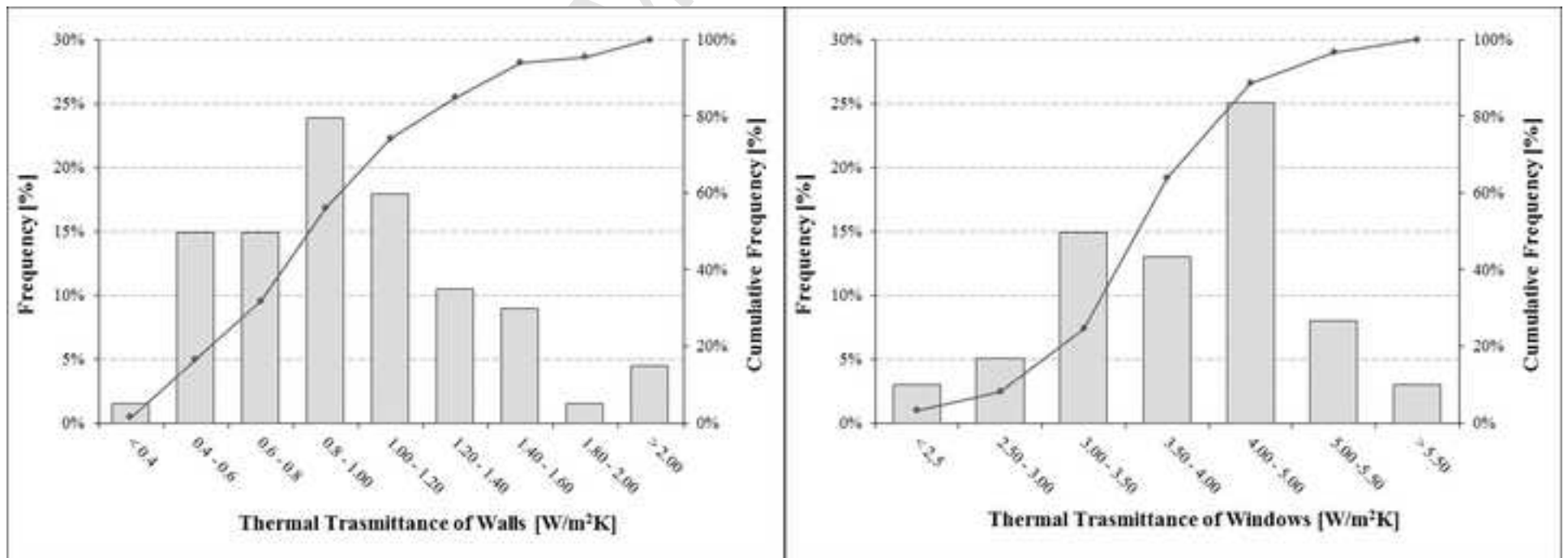


Figure 4
[Click here to download high resolution image](#)

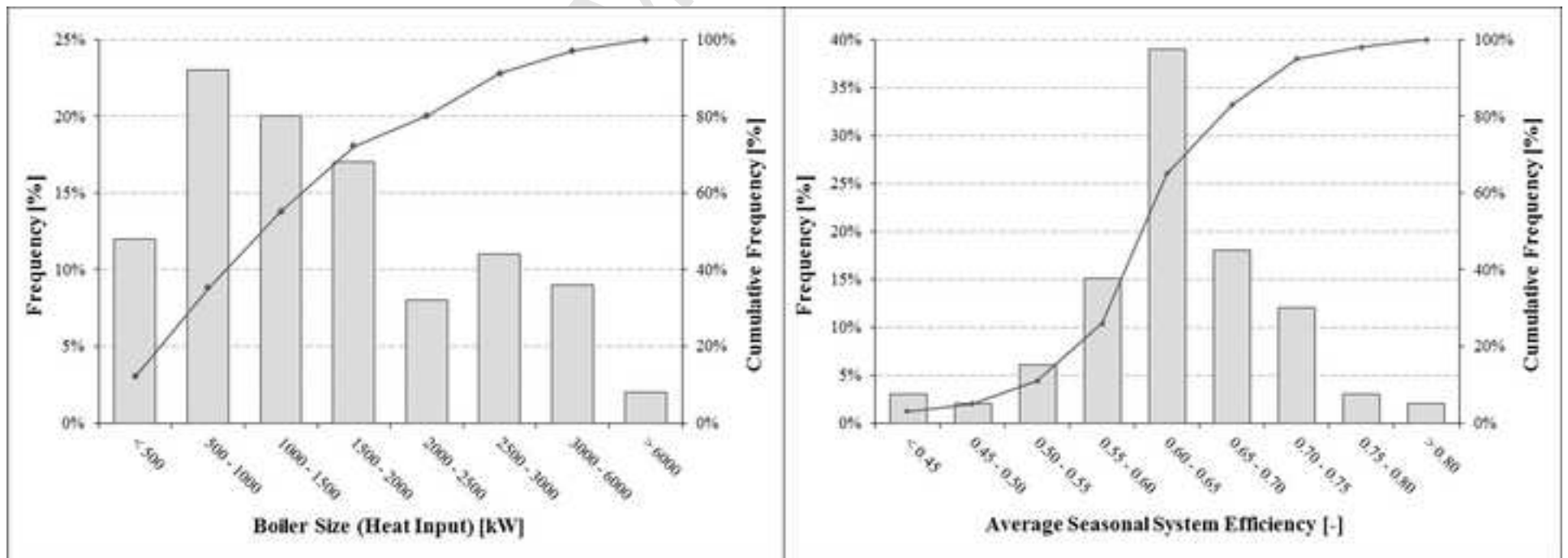


Figure 5
[Click here to download high resolution image](#)

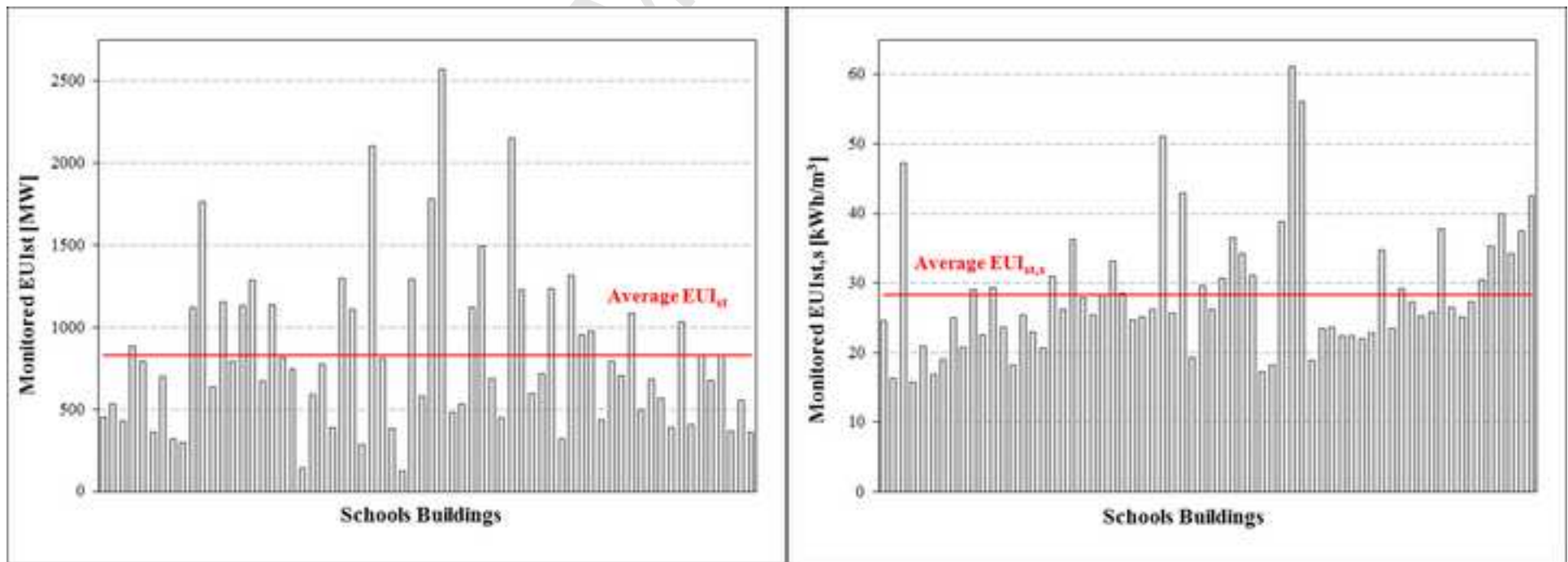


Figure 6
[Click here to download high resolution image](#)

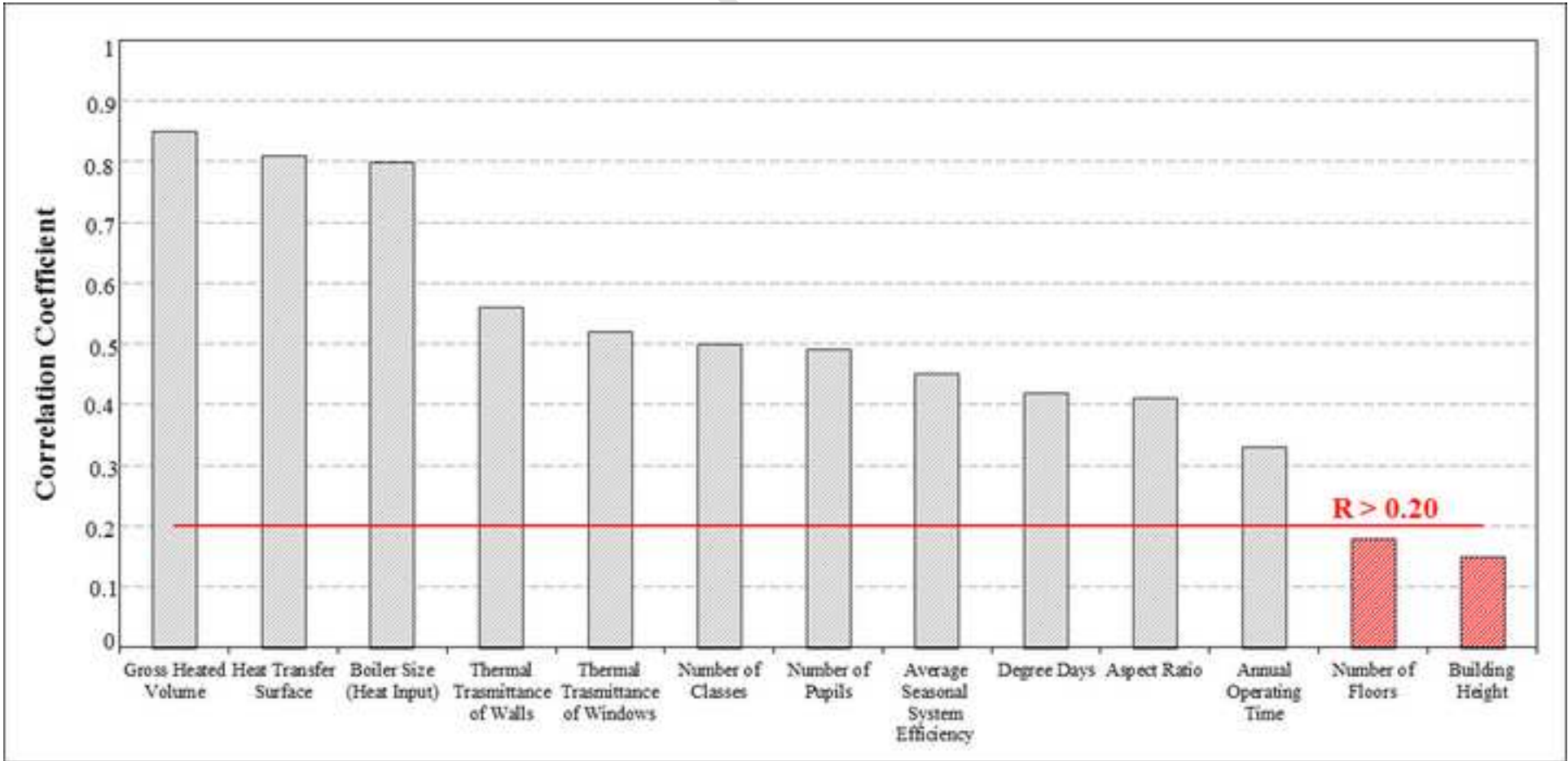


Figure 7
[Click here to download high resolution image](#)

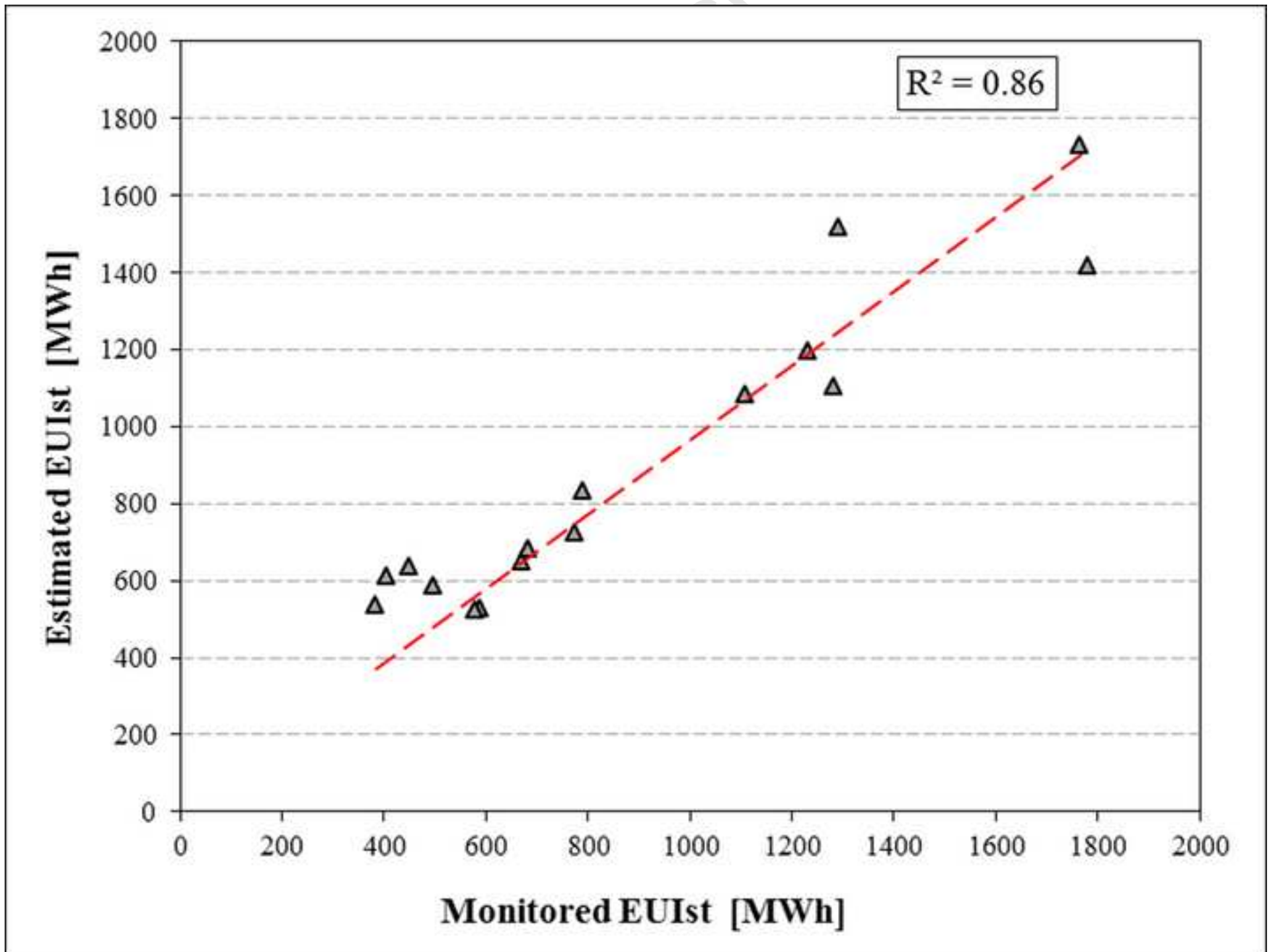


Figure 8
[Click here to download high resolution image](#)

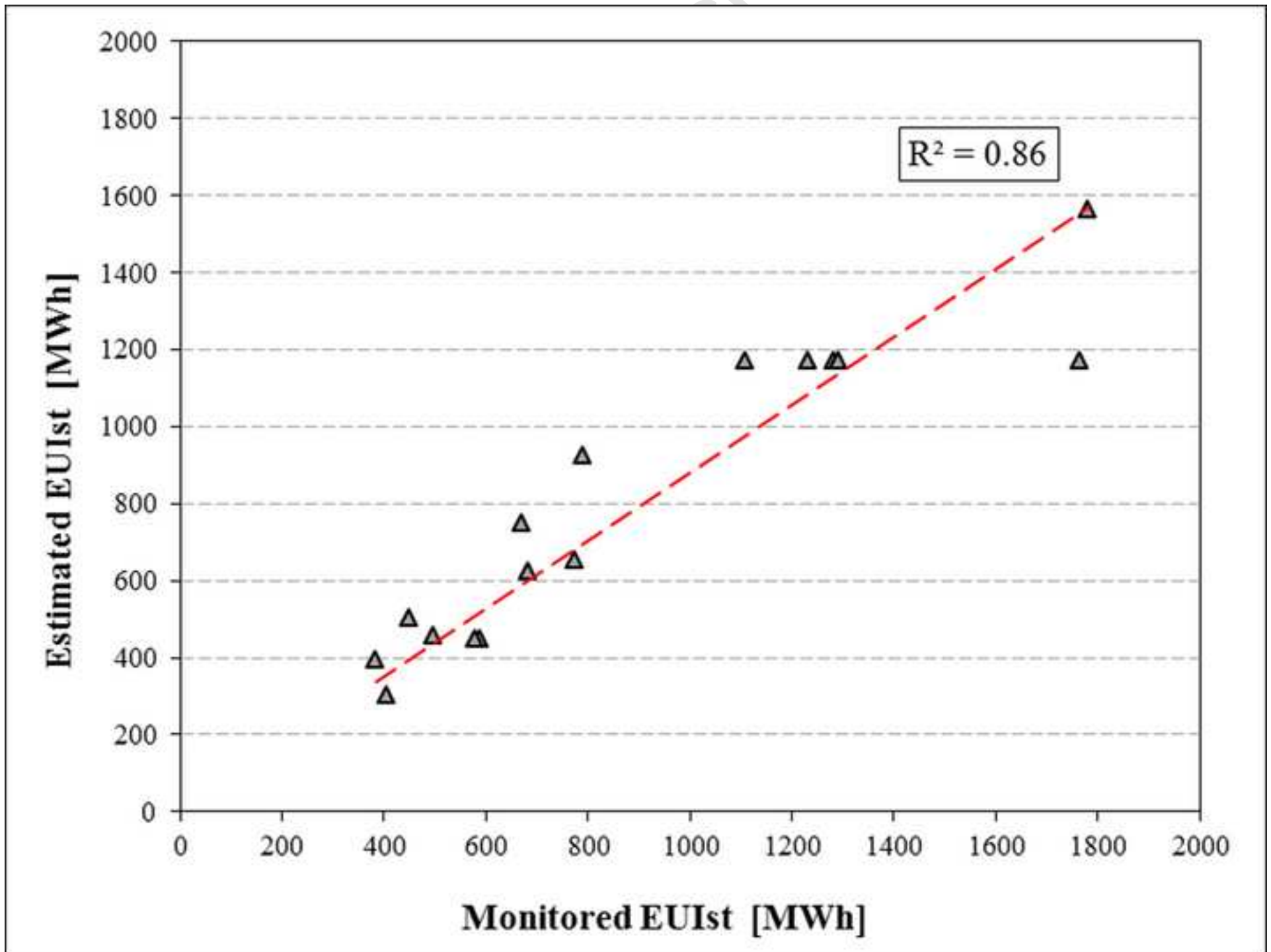


Figure 9
[Click here to download high resolution image](#)

