



(FP7 614100)

## **D5.2 Data Analysis and Forecast for Energy Consumption**

**29 July 2014 – Version 1.0**

**Published by the IMPReSS Consortium**

**Dissemination Level: Public**



**Project co-funded by the European Commission within the 7th Framework Programme and  
the Conselho Nacional de Desenvolvimento Científico e Tecnológico  
Objective ICT-2013.10.2 EU-Brazil research and development Cooperation**

## Document control page

**Document file:** D5.2 Data Analysis and Forecast for Energy Consumption.docx  
**Document version:** 1.0  
**Document owner:** Eduardo Souto (UFAM)

**Work package:** WP5 – Data Storage, Analysis & Decision Support  
**Task:** T5.2 Data Analysis and forecast for energy consumption  
**Deliverable type:** P

**Document status:** [X] approved by the document owner for internal review  
 [X] approved for submission to the EC

### Document history:

Version	Author(s)	Date	Summary of changes made
0.1	Eduardo Souto (UFAM)	13/07/2014	First Draft.
0.2	Thiago Rocha(UFAM)	15/07/2014	Second Draft.
0.3	Wesllen Sousa(UFAM)	17/07/2014	Third Draft.
0.4	Eulanda Santos(UFAM)	18/07/2014	Fixed Data Mining and Machine Learning concepts.
0.5	Francisco Ivan(UFAM)	19/07/2014	Fixed Optimization concepts.
0.6	Lucas Gomes (UFPE)	21/07/2014	Fixed typos and revised the document.
1.0	Eduardo Souto(UFAM)	29/07/2014	Final version submitted.

### Internal review history:

Reviewed by	Date	Summary of comments
Carlos Kamienski (UFABC)	26/07/2014	Accepted with minor corrections and comments
Ferry Pramudianto (FIT)	28/07/2014	Accepted with minor corrections and comments

### Legal Notice

The information in this document is subject to change without notice.

The Members of the IMPReSS Consortium make no warranty of any kind with regard to this document, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose. The Members of the IMPReSS Consortium shall not be held liable for errors contained herein or direct, indirect, special, incidental or consequential damages in connection with the furnishing, performance, or use of this material.

Possible inaccuracies of information are under the responsibility of the project. This report reflects solely the views of its authors. The European Commission is not liable for any use that may be made of the information contained therein.

## Index:

<b>1. Executive summary</b> .....	<b>4</b>
<b>2. Concepts</b> .....	<b>5</b>
2.1. Data Mining and Machine Learning .....	5
2.1.1. Classification.....	5
2.1.2. Regression.....	5
2.1.3. Cluster .....	5
2.1.4. Association .....	6
2.2. Optimization .....	6
2.2.1. Metaheuristic .....	6
2.2.2. Stochastic Algorithms.....	6
2.2.3. Evolutionary Algorithms.....	6
2.2.4. Swarm Algorithms .....	7
2.2. Statistics.....	7
<b>3. Architecture</b> .....	<b>8</b>
<b>4. Use Cases</b> .....	<b>16</b>
4.1. Scenario .....	16
4.1. First use case .....	17
4.2. Second use case .....	17
<b>5. References</b> .....	<b>22</b>

## 1. Executive summary

This deliverable focuses on showing how developers can use the IMPReSS platform to get insight from the massive amount of data acquired from smart objects connected to the IMPReSS Cloud. By the means of simulation scenarios, users will be capable of analyzing information provided by those smart objects, in order to extract knowledge from them.

Different statistic, optimization, data mining and machine learning tools may be used to help developers on extracting knowledge. First, this deliverable will show a general overview to highlight how IMPReSS will provide these strategies to the developers. Thus, it will focus on the statistics tools, since the optimization, data mining and machine learning tools will be discussed in the next deliverables.

In order to reach these goals, this deliverable is divided as follows. The first section discusses some concepts about statistics, optimization, data mining and machine learning algorithms. Then, sections two and three present the extended version of the API architecture for the Data Analytics module. Finally, section four shows some test case scenarios which were executed in order to lighten the developers. These scenarios are provided to help the developers to understand how they will use the Data Analytics module API, which is part of the IMPReSS API.

## 2. Concepts

This section discusses the main concepts employed in this deliverable. These concepts were used to understand and develop an extended version of the Data Analytics Application Programming Interface (API), which was first proposed in deliverable 8.2 [1]. First, data mining and machine learning concepts will be shown, then optimization and finally statistics.

### 2.1. Data Mining and Machine Learning

Extracting patterns and knowledge from large amounts of information is an important step for the solution of several tasks, such as remote sensing, financial forecasting, intrusion detection and fraud, information retrieval on the Internet, filtering spam emails, among others. Solutions proposed to solve such tasks are usually based on two inter-related fields: Machine Learning and Data Mining.

Data mining focuses on the discovery of (previously) unknown properties in the data, such as the analysis step of knowledge discovery in databases, while machine learning focuses on prediction, based on known properties learned from the training data. Machine Learning techniques can be divided into two groups: supervised and unsupervised. Supervised learning algorithms are most used to forecast a value and require a definition of a known outcome or target for the scenario. The objective of supervised learning is classification, for discrete target values, or regression, for continuous output. Unsupervised learning algorithms are usually used to recognize a model structure that is relevant in a set of data [2]. Clustering algorithms are examples of unsupervised learning methods. The Data Analytics module API will provide several classification, regression, association and cluster algorithms.

#### 2.1.1. Classification

Handwritten digits recognition, business modeling, and credit analysis are examples of classification problems. As a type of supervised learning, a classification algorithm builds a model that is used to assign labels to the unknown examples. The data used to generate the models may include energy consumption information such as cost and time to allow the prediction of classes of behaviors. The input or training data for a supervised learning algorithm requires the presence of attributes to represent each training sample, as well the classes assigned to the training samples [3].

#### 2.1.2. Regression

Regression algorithms are tools employed to predict a number as the output of a problem. Profit, sales, mortgage rates, house values, square footage, temperature or distance could all be predicted using regression techniques [4]. For example, a regression model could be used to predict the energy cost of a house based on number of rooms, number of equipment and other factors.

#### 2.1.3. Cluster

Clustering techniques try to identify natural clusters of data objects according to a similarity measure, for instance Euclidian distance. The members of the same cluster are more similar to each other than they are in comparison with members of other clusters. The goal of clustering analysis is to find high-quality clusters such that the inter-cluster similarity is low and the intra-cluster similarity is high [4].

#### **2.1.4. Association**

Association strategies are used to identify some attribute value conditions that occur frequently together in a given set of data. Association analysis is widely used in transaction data analysis for directed marketing, catalog design, and other business decision-making process. Traditionally, association is used for market basket data analysis such as 90% of the people who buy milk also buy bread [4].

### **2.2. Optimization**

Optimization consists on the effort of maximizing or minimizing an objective function by systematically choosing input values from within a desired range. And, thereafter, compute the value of the objective function. The main issue in optimization applications is to solve a problem finding a solution or a set of solutions with the most effective objective values (the lowest value in minimization problems or the highest value in maximization problems) under some given constraints [5]. The Data Analytics module API will provide a set of optimization algorithms to solve problems in energy consumption scenarios. Since these scenarios will typically involve a large initial amount of data, metaheuristics should be used due to the exponential complexity of an exhaustive search.

#### **2.2.1. Metaheuristic**

The objective of a heuristic is to produce a solution in a reasonable time that is good enough for solving the problem at hand. The chosen solution may not be the best among all possible solutions for the problem [6], i.e. there is no guarantee that the optimal solution will be found. However, metaheuristics traditionally employ strategies in order to create a process capable of escaping local minima and perform robust search in the solution space of a problem, thus seeking solutions very close to the global optimum solutions [7]. Greedy Search, Genetic Algorithms, Simulated Annealing and Particle Swarm Optimization are examples of the most popular metaheuristics. Evolutionary and Swarm optimization algorithms are described below.

#### **2.2.2. Stochastic Algorithms**

Stochastic algorithms are predominantly global optimization algorithms and metaheuristics that manage the application of local search procedures. This set of algorithms provides various different strategies to allow the generation of a variety of starting points in order to employ local search for refinement. This process is repeated leading to unexplored or potentially better areas of the search space [8]. Some examples of algorithms of this group are Random search, Adaptive Random Search, Stochastic Hill Climbing, Iterated Local Search, among others.

#### **2.2.3. Evolutionary Algorithms**

Evolutionary Algorithms belong to the Evolutionary Computation field of study and is concerned with computational methods inspired by the process and mechanisms of biological evolution. Similarly, as the process of evolution by means of natural selection proposed by Darwin. The mechanisms of evolution describe how evolution actually takes place through the modification and propagation of genetic material. Evolutionary Algorithms are concerned with investigating computational systems that resemble simplified versions of the processes and mechanisms of evolution toward achieving

the effects of these processes and mechanisms, namely the development of adaptive systems [8]. Therefore, Evolutionary Algorithms mimic nature's evolutionary principles to conduct its search towards an optimal solution, even though the optimal solution may not be achieved.

#### **2.2.4. Swarm Algorithms**

Swarm intelligence is the study of computational systems inspired by the 'collective intelligence'. Collective Intelligence emerges through the cooperation of large numbers of homogeneous agents in the environment. Examples include schools of fish, flocks of birds, and colonies of ants. Such intelligence is decentralized, self-organizing and distributed through out an environment. In nature, such systems are commonly used to solve problems such as e effective foraging for food, prey evading, or colony re-location [8]. Particle Swarm Optimization (PSO), for instance, simulates the behaviors of bird flocking or fish schooling. Each individual of the population is called particle. All particles have fitness values which are evaluated during the optimization process towards finding the best particle in the whole population.

### **2.2. Statistics**

Statistics is the study of the collection, organization, analysis, interpretation and presentation of data. When analyzing data, it is possible to use several statistics methods [9]. The reason to use statistics is to determine what information is reliable and which predictions can be trusted. The data stored has particular representation and relevant significance. Therefore, it is possible to apply statistic functions aiming to represent data from an energy consumption scenario, which may be used, for instance, to predict consumption or generate metrics to apply rules to solve a problem. The Data Analytics module API will provide several statistical methods.

### 3. Architecture

This section shows an overview of the Data Analytics module API. The main goal of this API is to provide to the developers tools related to data analytics, precisely statistics, optimization, data mining and machine learning algorithms. This API will be part of the IMPReSS Cloud and will be available through a web service. Figure 3.1 shows the Data Analytics module API Overview.

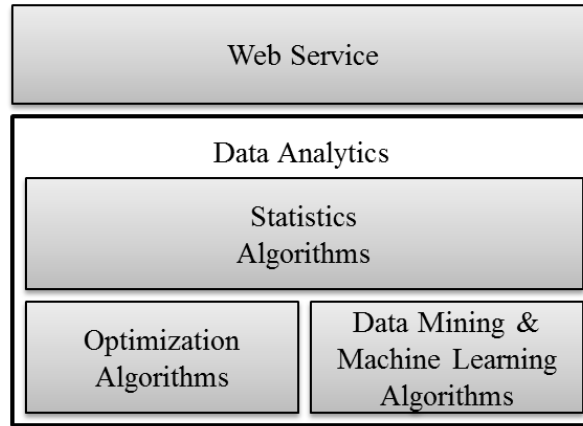


Figure 3.1. Data Analytics API Overview.

There is a web service, which implements the Data Analytics module API, in order to allow developers to have access to the Data Analytics module API. Deliverable 8.2 presents details about how the web services URLs should be formatted in order to guarantee consistency among IMPReSS' different modules. The second part of Figure 3.1 shows how the data analytics process is divided at the IMPReSS Platform. Data Analytics entails the process of examining large amounts of data to uncover hidden patterns, unknown correlations and other useful information that can be used to generate better decisions. With big data analytics, data scientists and other users can analyze huge volumes of data that conventional analytics and business intelligence solutions cannot do. This API offers four groups of algorithms to analyze raw data: statistics, optimization, data mining and machine learning algorithms. These groups are better depicted in Figure 3.2.



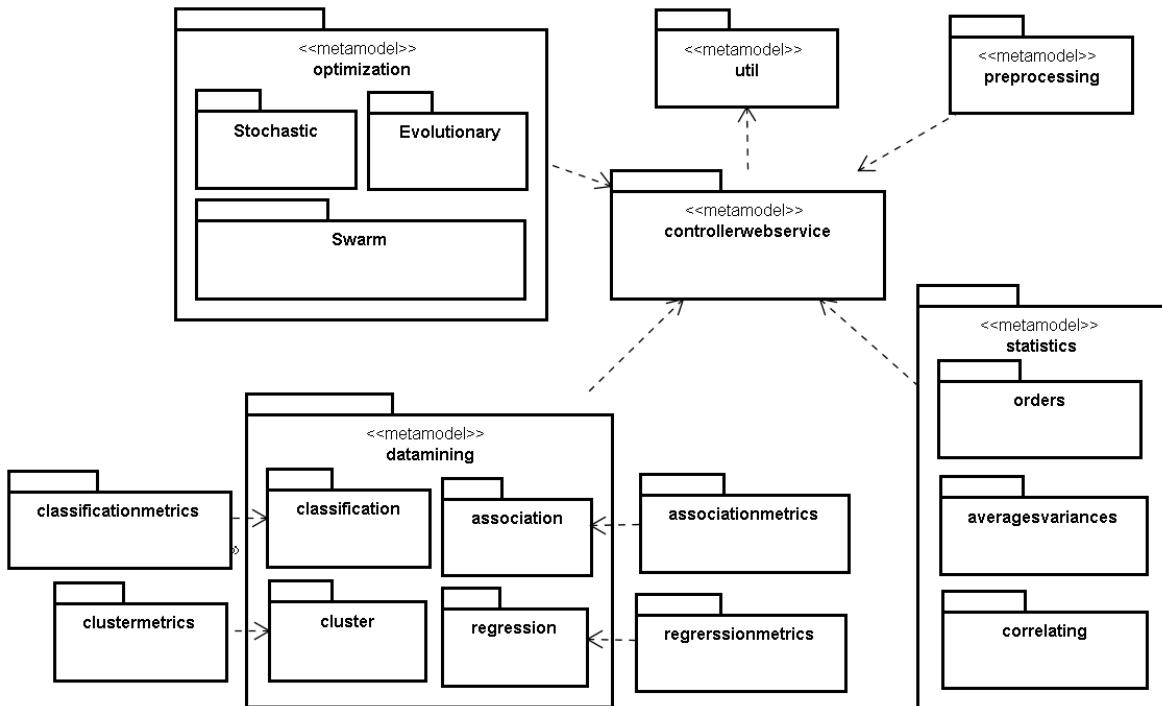


Figure 3.2. Data Analytics module API package structure.

The design of this API follows patterns from the Java Specification Request (JSR) [10] through the creation of packages. Furthermore, data mining algorithms were taken from the scikit-learn library [11], statistics algorithms from pymc library [12] and optimization algorithms from Clever Algorithms [8].

Before describing data mining, optimization and statistics, there are other packages that compose the Data Analytics module API. The first one is the **impress.util** package that consists of some utility classes which may be used to assist the other packages, for example, a class can be used to call a method to format data. The second package is the **impress.controllerwebservice** that has classes to expose the Data Analytics module API through a web service and makes it available to developers. Deliverable 8.2 has more details about the web service specification. The last package is the **impress.preprocessing** that consists of several classes to allow data preprocessing work before any data analytics algorithm is called. **Impress.preprocessing** algorithms are shown in Table 3.1. The preprocessing step may involve several strategies, for instance, when there are data with missing features, imputation methods should be employed to substitute missing values with meaningful estimates in order to allow machine learning algorithms to be used. Also, when data is available in different scales, normalization and standardization preprocessing should be used to allow the effective use of learning methods based on distance measures, such as Euclidian distance.

Table 3.1: Preprocessing algorithms.

Classes	Description
<b>Binarizer</b> ([threshold, copy])	Binarize data (set feature values to 0 or 1) according to a threshold
<b>Imputer</b> ([missing_values, ...])	Imputation methods for completing missing values
<b>KernelCenterer</b>	Center a kernel matrix
<b>LabelBinarizer</b> ([neg_label, ...])	Binarize labels in a one-vs-all fashion
<b>LabelEncoder</b>	Encode labels with value between 0 and n_classes-1.
<b>MinMaxScaler</b> ([feature_range, copy])	Standardizes features by scaling each feature to a given range.
<b>Normalizer</b> ([norm, copy])	Normalize samples individually to unit norm
<b>StandardScaler</b> ([copy, ...])	Standardize features by removing the mean and scaling to unit variance
<b>binarize</b> (X[, threshold, copy])	Boolean thresholding of array-like or scipy.sparse matrix

<b>label binarize</b> (y, classes[, ...])	Binarize labels in a one-vs-all fashion
<b>normalize</b> (X[, norm, axis, copy])	Normalize a dataset along any axis
<b>scale</b> (X[, axis, with_mean, ...])	Standardize a dataset along any axis

Figure 3.3 consists of several data mining and machine learning algorithms, as well as metrics that each group uses to evaluate their algorithms. These algorithms and metrics will be explained below:

- **impress.datamining.classification**: This sub-package defines classes that implements classification algorithms.
- **impress.datamining.association**: This sub-package defines classes that implements association algorithms.
- **impress.datamining.cluster**: This sub-package defines classes that implements clustering algorithms.
- **impress.datamining.regression**: This sub-package defines classes that implements regression algorithms.
- **impress.datamining.classification.metrics**: This sub-package defines metrics classes to evaluate data mining classification algorithms.
- **impress.datamining.association.metrics**: This sub-package defines metrics classes to evaluate association algorithms.
- **impress.datamining.cluster.metrics**: This sub-package defines metrics classes to evaluate data mining clustering algorithms.
- **impress.datamining.regression.metrics**: This sub-package defines metrics classes to evaluate data mining regression algorithms.

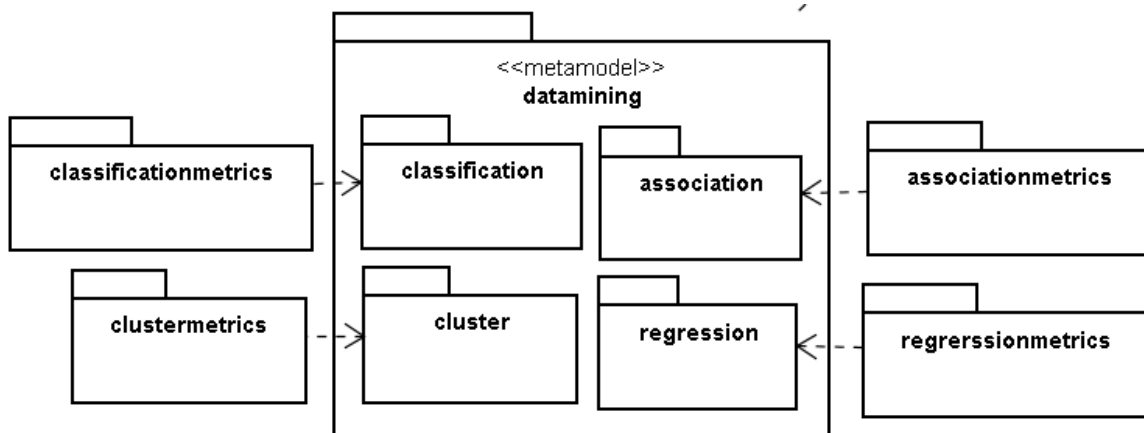


Figure 3.3. Data Mining package structure.

The algorithms of each package cited above are shown in the tables below:

Table 3.2: Classification algorithms.

Classes	Description
<b>RandomForestClassifier</b> ([...])	Ensemble of decision tree classifiers generated by random forest.
<b>RandomTreesEmbedding</b> ([...])	Ensemble of totally random trees.
<b>AdaBoostClassifier</b> (...[, criterion])	Ensemble of classifiers generated by the AdaBoost ensemble generation method.
<b>GradientBoostingClassifier</b> ([loss, ...])	Ensemble of classifiers generated by the Gradient Boosting method.
<b>GaussianHMM</b> ([n_components, ...])	The Hidden Markov Model classifier generated with Gaussian emissions

<b>GaussianNB</b>	The Naive Bayes Classifier generated with Gaussian functions
<b>NearestCentroid</b> ([metric, ...])	The prototype-based Nearest centroid classifier.
<b>SVC</b> ([C, kernel, degree, gamma, coef0, ...])	The Support Vector Machines classifier C-Support Vector Classification for nonlinear separated problems.
<b>LinearSVC</b> ([penalty, loss, dual, tol, C, ...])	Linear Support Vector Classification for linear separated problems.
<b>DecisionTreeClassifier</b> ([criterion, ...])	A decision tree classifier.

Table 3.3: Evaluation metrics for classification algorithms.

Classes	Description
<b>accuracy_score</b> (y_true, y_pred[, ...])	Accuracy classification score.
<b>average_precision_score</b> (y_true, y_score)	Compute average precision (AP) from prediction scores
<b>classification_report</b> (y_true, y_pred)	Build a text report showing the main classification metrics
<b>confusion_matrix</b> (y_true, y_pred[, ...])	Compute confusion matrix to evaluate the accuracy of a classification
<b>jaccard_similarity_score</b> (y_true, y_pred)	Jaccard similarity coefficient score
<b>log_loss</b> (y_true, y_pred[, eps, ...])	Log loss, aka logistic loss or cross-entropy loss.
<b>precision_recall_curve</b> (y_true, ...)	Compute precision-recall pairs for different probability thresholds
<b>precision_score</b> (y_true, y_pred[, ...])	Compute the precision
<b>recall_score</b> (y_true, y_pred[, ...])	Compute the recall
<b>roc_auc_score</b> (y_true, y_score)	Compute Area Under the Curve (AUC) from prediction scores

Table 3.4: Clustering algorithms.

Classes	Description
<b>DBSCAN</b> ([eps, min_samples, metric, ...])	Perform DBSCAN clustering from vector array or distance matrix.
<b>KMeans</b> ([n_clusters, init, n_init, ...])	The simple K-Means clustering algorithm
<b>MiniBatchKMeans</b> ([n_clusters, init, ...])	Mini-Batch K-Means clustering algorithm
<b>MeanShift</b> ([bandwidth, seeds, ...])	MeanShift clustering algorithm
<b>Ward</b> ([n_clusters, memory, ...])	Ward hierarchical clustering: constructs a tree and cuts it.
<b>NearestNeighbors</b> ([n_neighbors, ...])	Unsupervised learner for implementing neighbor searches.

Table 3.5: Evaluation metrics for Clustering algorithms.

Classes	Description
<b>adjusted_mutual_info_score</b> (...)	Adjusted Mutual Information between two clusters
<b>adjusted_rand_score</b> (labels_true, ...)	Rand index adjusted for chance
<b>completeness_score</b> (labels_true, ...)	Completeness metric of a cluster labeling given a ground truth
<b>homogeneity_score</b> (labels_true, ...)	Homogeneity metric of a cluster labeling given a ground truth
<b>mutual_info_score</b> (labels_true, ...)	Mutual Information between two clusters
<b>normalized_mutual_info_score</b> (...)	Normalized Mutual Information between two clusters
<b>consensus_score</b> (a, b[, similarity])	The similarity of two sets of biclusters.

Table 3.6: Regression algorithms.

Classes	Description
<b>RandomForestRegressor</b> ([...])	Ensemble of random forest regressor.
<b>AdaBoostRegressor</b> (...[, criterion, ...])	Ensemble based on AdaBoost regressor.
<b>GradientBoostingRegressor</b> ([loss, ...])	Gradient Boosting for regression.
<b>IsotonicRegression</b> ([y_min, y_max, ...])	Isotonic regression model.
<b>LinearRegression</b> ([...])	Ordinary least squares Linear Regression.
<b>LogisticRegression</b> ([penalty, ...])	Logistic Regression classifier.
<b>RandomizedLogisticRegression</b> ([...])	Randomized Logistic Regression
<b>KNeighborsRegressor</b> ([n_neighbors, ...])	Regression based on k-nearest neighbors.
<b>RadiusNeighborsRegressor</b> ([radius, ...])	Regression based on neighbors within a fixed radius.
<b>SVR</b> ([kernel, degree, gamma, coef0, tol, ...])	epsilon-Support Vector Regression.

Table 3.7: Evaluation metrics for Regression algorithms.

Classes	Description
<b>explained_variance_score</b> (y_true, y_pred)	Explained variance regression score function
<b>mean_absolute_error</b> (y_true, y_pred)	Mean absolute error regression loss
<b>mean_squared_error</b> (y_true, y_pred)	Mean squared error regression loss
<b>r2_score</b> (y_true, y_pred)	R <sup>2</sup> (coefficient of determination) regression score function.

Package **impress.optimization** consists of several optimization algorithms, such as Tabu Search, Evolutionary Algorithms and Particle Swarm Optimization, as depicted in Figure 3.4 and explained below.

- **impress.optimization.stochastic:** This sub-package defines classes that implements stochastic optimization algorithms.
- **impress.optimization.evolutionary:** This sub-package defines classes that implements evolutionary optimization algorithms.
- **impress.optimization.swarm:** This sub-package defines classes that implements swarm optimization algorithms

The algorithms of each of these packages are shown in the tables below:

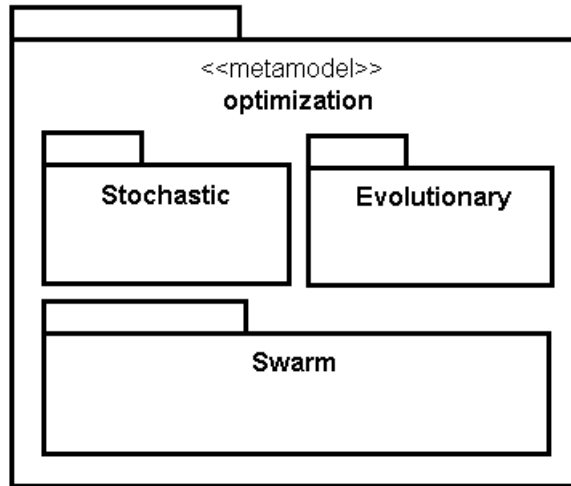


Figure 3.4. Data analytics module optimization package structure.

Table 3.8: Stochastic algorithms.

Classes	Description
<b>tabusearch</b> (<parameters>)	Tabu Search is a Global Optimization algorithm and a Metaheuristic or Meta-strategy for controlling an embedded heuristic technique.
<b>randomsearch</b> (<parameters>)	Random search belongs to the elds of Stochastic Optimization and Global Optimization.
<b>scattersearch</b> (<parameters>)	Scatter search is a Metaheuristic and a Global Optimization algorithm.

Table 3.9: Evolutionary algorithms.

Classes	Description
<b>geneticalgorithm</b> (<parameters>)	The Genetic Algorithm is an Adaptive Strategy and a Global Optimization technique.
<b>nsga</b> (<parameters>)	The Non-dominated Sorting Genetic Algorithm is a Multiple Objective Optimization (MOO) algorithm and is an instance of an Evolutionary Algorithm from the field of Evolutionary Computation.
<b>spea</b> (<parameters>)	Strength Pareto Evolutionary Algorithm is a Multiple Objective Optimization (MOO) algorithm and an Evolutionary Algorithm from the field of Evolutionary Computation.

Table 3.10: Swarm algorithms.

Classes	Description
<b>pso</b> (<parameters>)	Particle Swarm Optimization belongs to the field of Swarm Intelligence and Collective Intelligence and is a sub-field of Computational Intelligence.
<b>beesalgorithm</b> (<parameters>)	The Bees Algorithm beings to Bee Inspired Algorithms and the field of Swarm Intelligence, and more broadly the elds of Computational Intelligence and Metaheuristics.
<b>antcolony</b> system(<parameters>)	The Ant Colony System algorithm is an example of an Ant Colony Optimization method from the field of Swarm Intelligence, Metaheuristics and Computational Intelligence.

Package **impress.statistics** consists of several statistics algorithms, such as orders, average, variance and correlating, as depicted in Figure 3.5 and explained below.

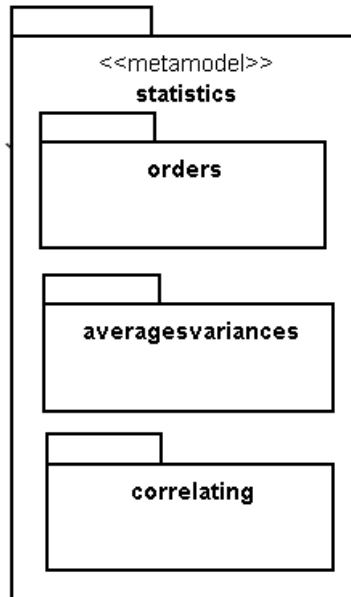


Figure 3.5. IMPRESS statistic package structure.

- **impress.statistics.orders:** This sub-package defines classes that implement the extraction of basics information, i.e. minimum and maximum.
- **impress.statistics.averagevariances:** This sub-package defines classes that implement the extraction of average and variances of the data.
- **impress.statistics.correlating:** This sub-package defines classes that implement the extraction of correlation coefficients.

The algorithms of each package are shown in the tables bellow:

Table 3.11: Order algorithms.

Classes	Description
<b>amin</b> (a[, axis, out, keepdims])	Returns the minimum of an array or minimum along an axis.
<b>amax</b> (a[, axis, out, keepdims])	Returns the maximum of an array or maximum along an axis.
<b>nanmin</b> (a[, axis, out, keepdims])	Returns minimum of an array or minimum along an axis, ignoring any NaNs.
<b>nanmax</b> (a[, axis, out, keepdims])	Returns the maximum of an array or maximum along an axis, ignoring any
<b>ptp</b> (a[, axis, out])	Range of values (maximum - minimum) along an axis.
<b>percentile</b> (a, q[, axis, out, overwrite_input])	Compute the qth percentile of the data along the specified axis.

Table 3.12: Average/Variance algorithms.

Classes	Description
<b>median</b> (a[, axis, out, overwrite_input])	Computes the median along the specified axis.
<b>average</b> (a[, axis, weights, returned])	Computes the weighted average along the specified axis.
<b>mean</b> (a[, axis, dtype, out, keepdims])	Computes the arithmetic mean along the specified axis.

<b>std</b> (a[, axis, dtype, out, ddof, keepdims])	Computes the standard deviation along the specified axis.
<b>var</b> (a[, axis, dtype, out, ddof, keepdims])	Computes the variance along the specified axis.
<b>nanmean</b> (a[, axis, dtype, out, keepdims])	Computes the arithmetic mean along the specified axis, ignoring NaNs.
<b>nanstd</b> (a[, axis, dtype, out, ddof, keepdims])	Computes the standard deviation along the specified axis, while
<b>nanvar</b> (a[, axis, dtype, out, ddof, keepdims])	Computes the variance along the specified axis, while ignoring NaNs.

Table 3.13: Correlating algorithms.

Classes	Description
<b>corrcoef</b> (x[, y, rowvar, bias, ddof])	Returns correlation coefficients.
<b>correlate</b> (a, v[, mode, old_behavior])	Cross-correlation of two 1-dimensional sequences.
<b>cov</b> (m[, y, rowvar, bias, ddof])	Estimates a covariance matrix, given data.

The next section focuses on some scenarios to show to the developers how they can use the Data Analytics module API to use any of the aforementioned algorithms. These scenarios, though, are focused on statistics algorithms. Optimization, data mining and machine learning algorithms will be worked with details in the next deliverables of WP5.

## 4. Use Cases

The potential applications of Data Analytics module API are wide-ranging. This section covers two use cases to show the usage of the statistical capabilities of the aforementioned API. A scenario with a database populated with real data based on a study from Massachusetts Institute of Technology (MIT) [13] was chosen to address the use of some data analytics features focused at the statistical algorithms such as standard deviation, mean, maximum, minimum, among others.

### 4.1. Scenario

The scenario is composed by 77 sensors installed in an apartment to collect data about human activity everyday. For the purpose of this scenario, only data related to appliance energy consumption was collected. Figure 4.1 has a preview of the dataset created from the data collected.

id integer	name character var	model character var	consumption numeric	priority integer	start_time bigint	end_time bigint	duration bigint	locations character var	day character var
35	Dryer Cloth	BSR10 10 kg	4800	4	31483	3895	-27588	Dinning room	monday
26	Dishwasher	Brastemp At	780	3	31483	3895	-27588	Dinning room	thursday
36	Toaster	Daily Ri25	800	3	31483	3895	-27588	Bathroom	saturday
25	washes cloth	Brastemp At	300	1	31483	3895	-27588	Living room	wednesday
35	Dryer Cloth	BSR10 10 kg	4800	4	31483	3895	-27588	Bathroom	sunday
20	Cooker	Electrolux	1180	4	31483	3895	-27588	Living room	monday
36	Toaster	Daily Ri25	800	3	31483	3895	-27588	Bathroom	friday
28	Sewing mach	Singer Brill	70	1	31483	3895	-27588	Dinning room	tuesday
35	Dryer Cloth	BSR10 10 kg	4800	4	31483	3895	-27588	Dinning room	saturday
5	Air conditi	Split BBV0	822	3	31483	3895	-27588	Kitchen	monday
36	Toaster	Daily Ri25	800	3	31483	3895	-27588	Bathroom	thursday
24	Table lamp	Boyu	30	1	31483	3895	-27588	Kitchen	sunday
35	Dryer Cloth	BSR10 10 kg	4800	4	31483	3895	-27588	Dinning room	thursday
20	Cooker	Electrolux	1180	4	31483	3895	-27588	Living room	sunday
36	Toaster	Daily Ri25	800	3	31483	3895	-27588	Bathroom	wednesday
17	Electric fr	Britania Pr	90	1	31483	3895	-27588	Living room	saturday
35	Dryer Cloth	BSR10 10 kg	4800	4	31483	3895	-27588	Dinning room	friday

Figure 4.1. Dataset from appliance energy consumption.

The attributes from the dataset are explained below:

- **id:** Appliance identification id.
- **name:** Appliance classification by name.
- **model:** Appliance model identification.
- **Consumption\_on:** Consumption about each appliance.
- **priority:** Priority definition for power consumption.
- **start\_time:** Time at which the activity started in seconds.
- **end\_time:** Time at which the activity ended in seconds.
- **duration:** Difference between start\_time and end\_time.
- **locations:** Appliance location.
- **day:** Number between 1 and 7 specifying the day of the week (1=monday, 7=sunday).



### 4.1. First use case

The first use case predicts the energy mean consumption per room for the next day, which is useful to know the average consumption and get a sense of how much this consumption is going to cost. To calculate this, the "mean" function from the Data Analytics module has to be called through the URL below.

*http://<domain>/statistics/averagesvariances/mean/listnumber?values=[a,b,c]*

This URL creates a request to the IMPRESS Platform to execute the mean function. Their parameters are explained in Table 4.1.

Table 4.1: Mean function parameters

<b>Parameters:</b>	<b>values</b> : <i>array_like</i> Array containing numbers whose mean is desired. If <i>a</i> is not an array, a conversion is attempted.
<b>Returns :</b>	<b>m</b> : <i>ndarray, see dtype parameter above</i>

After the call the developer gets a response from the Data Analytics module with the mean energy consumption per room for the next day. Figure 4.3 shows an example of response in JSON format, this response is not final and can change with time.

```

Response
Content-
Type: application/json
Date: <date>
Server: <server>
{mean:
{
result: X
}
}
    
```

Figure 4.3. JSON return from an API call..

### 4.2. Second use case

The second use case reports the maximum and minimum consumption of each room, which is useful to know which place is consuming more energy. To calculate this, the "min" and "max" functions from the Data Analytics module have to be called, an example for the min function is depicted through the URL below.

*http://<domain>/statistics/averagesvariances/amin/listnumber?values=[a,b,c]*

Table 4.2: Min function parameters

<b>Parameters :</b>	<b>values</b> : <i>array_like</i> Input data.
<b>Returns :</b>	<b>amin</b> : <i>ndarray or scalar</i>

After the call the developer gets a response from the Data Analytics module with the min consumption of each room. Figure 4.4 shows an example of response in JSON format, this response is not final and can change with time.

Response
<pre>Content-Type: application/json Date: &lt;date&gt; Server: &lt;server&gt; {min: { result: X } }</pre>

Figure 4.4. JSON return from an API call.

All the requests that can be made at the IMPRESS Platform are shown at Table 4.3, this list is not final and can change with time.

Table 4.3: Data Analytics module algorithms endpoints

Algorithms	URL
<b>Data Mining Cluster</b>	
<b>DBSCAN</b> ([eps, min_samples, metric, ...])	/datamining/cluster/ <b>DBSCAN</b> / <b>&lt;Parameters&gt;</b>
<b>KMeans</b> ([n_clusters, init, n_init, ...])	/datamining/cluster/ <b>KMeans</b> / <b>&lt;Parameters&gt;</b>
<b>MiniBatchKMeans</b> ([n_clusters, init, ...])	/datamining/cluster/ <b>MiniBatchKMeans</b> / <b>&lt;Parameters&gt;</b>
<b>MeanShift</b> ([bandwidth, seeds, ...])	/datamining/cluster/ <b>MeanShift</b> / <b>&lt;Parameters&gt;</b>
<b>Ward</b> ([n_clusters, memory, ...])	/datamining/cluster/ <b>Ward</b> / <b>&lt;Parameters&gt;</b>
<b>NearestNeighbors</b> ([n_neighbors, ...])	/datamining/cluster/ <b>NearestNeighbors</b> / <b>&lt;Parameters&gt;</b>
<b>Data Mining Cluster Metrics</b>	
<b>adjusted_mutual_info_score</b> (...)	/datamining/cluster/metrics/ <b>adjusted_mutual_info_score</b> / <b>&lt;Parameters&gt;</b>
<b>adjusted_rand_score</b> (labels_true, ...)	/datamining/cluster/metrics/ <b>adjusted_rand_score</b> / <b>&lt;Parameters&gt;</b>
<b>completeness_score</b> (labels_true, ...)	/datamining/cluster/metrics/ <b>completeness_score</b> / <b>&lt;Parameters&gt;</b>
<b>homogeneity_score</b> (labels_true, ...)	/datamining/cluster/metrics/ <b>homogeneity_score</b> / <b>&lt;Parameters&gt;</b>
<b>mutual_info_score</b> (labels_true, ...)	/datamining/cluster/metrics/ <b>mutual_info_score</b> / <b>&lt;Parameters&gt;</b>
<b>normalized_mutual_info_score</b> (...)	/datamining/cluster/metrics/ <b>normalized_mutual_info_score</b> / <b>&lt;Parameters&gt;</b>
<b>consensus_score</b> (a, b[, similarity])	/datamining/cluster/metrics/ <b>consensus_score</b> / <b>&lt;Parameters&gt;</b>
<b>Data Mining Classification</b>	
<b>RandomForestClassifier</b> ([...])	datamining/classification/ <b>RandomForestClassifier</b> / <b>&lt;Parameters&gt;</b>

<a href="#"><u>RandomTreesEmbedding</u></a> ([...])	datamining/classification/ <a href="#"><u>RandomTreesEmbedding</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>ExtraTreesClassifier</u></a> ([...])	datamining/classification/ <a href="#"><u>ExtraTreesClassifier</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>AdaBoostClassifier</u></a> ([...], criterion)	datamining/classification/ <a href="#"><u>AdaBoostClassifier</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>GradientBoostingClassifier</u></a> ([loss, ...])	datamining/classification/ <a href="#"><u>GradientBoostingClassifier</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>GaussianHMM</u></a> ([n_components, ...])	datamining/classification/ <a href="#"><u>GaussianHMM</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>GaussianNB</u></a>	datamining/classification/ <a href="#"><u>GaussianNB</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>NearestCentroid</u></a> ([metric, ...])	datamining/classification/ <a href="#"><u>NearestCentroid</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>SVC</u></a> ([C, kernel, degree, gamma, coef0, ..])	datamining/classification/ <a href="#"><u>SVC</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>LinearSVC</u></a> ([penalty, loss, dual, tol, C, ...])	datamining/classification/ <a href="#"><u>LinearSVC</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>DecisionTreeClassifier</u></a> ([criterion, ...])	datamining/classification/ <a href="#"><u>DecisionTreeClassifier</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<b>Data Mining Classification Metrics</b>	
<a href="#"><u>accuracy_score</u></a> (y_true, y_pred, ...)	datamining/classification/metrics/ <a href="#"><u>accuracy_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>average_precision_score</u></a> (y_true, y_score)	datamining/classification/metrics/ <a href="#"><u>average_precision_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>classification_report</u></a> (y_true, y_pred)	datamining/classification/metrics/ <a href="#"><u>classification_report</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>confusion_matrix</u></a> (y_true, y_pred, ...)	datamining/classification/metrics/ <a href="#"><u>confusion_matrix</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>jaccard_similarity_score</u></a> (y_true, y_pred)	datamining/classification/metrics/ <a href="#"><u>jaccard_similarity_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>precision_score</u></a> (y_true, y_pred, ...)	datamining/classification/metrics/ <a href="#"><u>precision_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>recall_score</u></a> (y_true, y_pred, ...)	datamining/classification/metrics/ <a href="#"><u>recall_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>roc_auc_score</u></a> (y_true, y_score)	datamining/classification/metrics/ <a href="#"><u>roc_auc_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<b>Data Mining Regression</b>	
<a href="#"><u>RandomForestRegressor</u></a> ([...])	/datamining/regression/ <a href="#"><u>RandomForestRegressor</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>AdaBoostRegressor</u></a> ([...], criterion, ...)	/datamining/regression/ <a href="#"><u>AdaBoostRegressor</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>GradientBoostingRegressor</u></a> ([loss, ...])	/datamining/regression/ <a href="#"><u>GradientBoostingRegressor</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>IsotonicRegression</u></a> ([y_min, y_max, ...])	/datamining/regression/ <a href="#"><u>IsotonicRegression</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>LinearRegression</u></a> ([...])	/datamining/regression/ <a href="#"><u>LinearRegression</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>LogisticRegression</u></a> ([penalty, ...])	/datamining/regression/ <a href="#"><u>LogisticRegression</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>RandomizedLogisticRegression</u></a> ([...])	/datamining/regression/ <a href="#"><u>RandomizedLogisticRegression</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>KNeighborsRegressor</u></a> ([n_neighbors, ..])	/datamining/regression/ <a href="#"><u>KNeighborsRegressor</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>RadiusNeighborsRegressor</u></a> ([radius, ...])	/datamining/regression/ <a href="#"><u>RadiusNeighborsRegressor</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>cross_decomposition.PLSRegression</u></a> ([...])	/datamining/regression/ <a href="#"><u>cross_decomposition.PLSRegression</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<a href="#"><u>SVR</u></a> ([kernel, degree, gamma, coef0, tol, ...])	/datamining/regression/ <a href="#"><u>SVR</u></a> / <a href="#"><u>&lt;Parameters</u></a> >
<b>Data Mining Regression Metrics</b>	
<a href="#"><u>explained_variance_score</u></a> (y_true, y_p)	/datamining/regression/metrics/ <a href="#"><u>explained_variance_score</u></a> / <a href="#"><u>&lt;Parameters</u></a> >

red)	meters>
<b>mean_absolute_error</b> (y_true, y_pred)	/datamining/regression/metrics/ <b>mean_absolute_error</b> /<<Parameters>
<b>mean_squared_error</b> (y_true, y_pred)	/datamining/regression/metrics/ <b>mean_squared_error</b> /<<Parameters>
<b>r2_score</b> (y_true, y_pred)	/datamining/regression/metrics/ <b>r2_score</b> /<<Parameters>
<b>Statistics Orders</b>	
<b>amin</b> (a[, axis, out, keepdims])	/statistics/orders/ <b>amin</b> /<<Parameters>
<b>amax</b> (a[, axis, out, keepdims])	/statistics/orders/ <b>amax</b> /<<Parameters>
<b>nanmin</b> (a[, axis, out, keepdims])	/statistics/orders/ <b>nanmin</b> /<<Parameters>
<b>nanmax</b> (a[, axis, out, keepdims])	/statistics/orders/ <b>nanmax</b> /<<Parameters>
<b>ptp</b> (a[, axis, out])	/statistics/orders/ <b>ptp</b> /<<Parameters>
<b>percentile</b> (a, q[, axis, out, overwrite_input])	/statistics/orders/ <b>percentile</b> /<<Parameters>
<b>Statistics Averages and Variances</b>	
<b>median</b> (a[, axis, out, overwrite_input])	/statistics/averagesvariances/ <b>median</b> /<<Parameters>
<b>average</b> (a[, axis, weights, returned])	/statistics/averagesvariances/ <b>average</b> /<<Parameters>
<b>mean</b> (a[, axis, dtype, out, keepdims])	/statistics/averagesvariances/ <b>mean</b> /<<Parameters>
<b>std</b> (a[, axis, dtype, out, ddof, keepdims])	/statistics/averagesvariances/ <b>std</b> /<<Parameters>
<b>var</b> (a[, axis, dtype, out, ddof, keepdims])	/statistics/averagesvariances/ <b>var</b> /<<Parameters>
<b>nanmean</b> (a[, axis, dtype, out, keepdims])	/statistics/averagesvariances/ <b>nanmean</b> /<<Parameters>
<b>nanstd</b> (a[, axis, dtype, out, ddof, keepdims])	/statistics/averagesvariances/ <b>nanstd</b> /<<Parameters>
<b>nanvar</b> (a[, axis, dtype, out, ddof, keepdims])	/statistics/averagesvariances/ <b>nanvar</b> /<<Parameters>
<b>Statistics Correlating</b>	
<b>corrcoef</b> (x[, y, rowvar, bias, ddof])	/statistics/correlation/ <b>corrcoef</b> /<<Parameters>
<b>correlate</b> (a, v[, mode, old_behavior])	/statistics/correlation/ <b>correlate</b> /<<Parameters>
<b>cov</b> (m[, y, rowvar, bias, ddof])	/statistics/correlation/ <b>cov</b> /<<Parameters>
<b>Optimization Stochastic Algorithms</b>	
<b>tabusearch</b> (<parameters>)	/optimization/stochastic/ <b>tabusearch</b> /<Parameters>
<b>randomsearch</b> (<parameters>)	/optimization/stochastic/ <b>randomsearch</b> /<Parameters>
<b>scattersearch</b> (<parameters>)	/optimization/stochastic/ <b>scattersearch</b> /<Parameters>
<b>Optimization Evolutionary Algorithms</b>	
<b>geneticalgorithm</b> (<parameters>)	/optimization/evolutionary/ <b>geneticalgorithm</b> /<Parameters>
<b>nsga</b> (<parameters>)	/optimization/evolutionary/ <b>nsga</b> /<Parameters>
<b>spea</b> (<parameters>)	/optimization/evolutionary/ <b>spea</b> /<Parameters>

<b>Optimization Swarm Algorithms</b>	
<b>pso</b> (<parameters>)	/optmization/swarm//<Parameters>
<b>beesalgorithm</b> (<parameters>)	/optmization/swarm//<Parameters>
<b>antcolony</b> system(<parameters>)	/optmization/swarm//<Parameters>
<b>Pre processing</b>	
<b>Binarizer</b> ([threshold, copy])	/preprocessing/ <b>Binarizer</b> / <b>&lt;Parameters&gt;</b>
<b>Imputer</b> ([missing_values, ...])	/preprocessing/ <b>Imputer</b> / <b>&lt;Parameters&gt;</b>
<b>KernelCenterer</b>	/preprocessing/ <b>KernelCenterer</b> / <b>&lt;Parameters&gt;</b>
<b>LabelBinarizer</b> ([neg_label, ...])	/preprocessing/ <b>LabelBinarizer</b> / <b>&lt;Parameters&gt;</b>
<b>LabelEncoder</b>	/preprocessing/ <b>LabelEncoder</b> / <b>&lt;Parameters&gt;</b>
<b>MinMaxScaler</b> ([feature_range, copy])	/preprocessing/ <b>MinMaxScaler</b> / <b>&lt;Parameters&gt;</b>
<b>Normalizer</b> ([norm, copy])	/preprocessing/ <b>Normalizer</b> / <b>&lt;Parameters&gt;</b>
<b>StandardScaler</b> ([copy, ...])	/preprocessing/ <b>StandardScaler</b> / <b>&lt;Parameters&gt;</b>
<b>binarize</b> (X[, threshold, copy])	/preprocessing/ <b>binarize</b> / <b>&lt;Parameters&gt;</b>
<b>label binarize</b> (y, classes[, ...])	/preprocessing/ <b>label binarize</b> / <b>&lt;Parameters&gt;</b>
<b>normalize</b> (X[, norm, axis, copy])	/preprocessing/ <b>normalize</b> / <b>&lt;Parameters&gt;</b>
<b>scale</b> (X[, axis, with_mean, ...])	/preprocessing/ <b>scale</b> / <b>&lt;Parameters&gt;</b>

## 5. References

- [1] IMPReSS 8.2 deliverable "Application Architecture for Energy Management", 14 May 2014.
- [2] Ian H. Witten, Eibe Frank, Mark A. Hall. Data Mining: Practical Machine Learning Tools and Techniques. Third Edition. 2011.
- [3] Oracle® Database Concepts 10g Release 1 (10.1). Available in: [http://docs.oracle.com/pdf/B10698\\_01.pdf](http://docs.oracle.com/pdf/B10698_01.pdf). Date Access: 07/11/2014.
- [4] Oracle® Database Concepts 11g Release 1 (11.1). Available in: [http://docs.oracle.com/cd/B28359\\_01/server.111/b28318.pdf](http://docs.oracle.com/cd/B28359_01/server.111/b28318.pdf). Date Access: 07/11/2014.
- [5] Business Dictionary. Available in: <http://www.businessdictionary.com/definition/optimization.html>. Date Access: 07/11/2014.
- [6] Russell, Stuart J.; Norvig, Peter (2003), Artificial Intelligence: A Modern Approach (2nd ed.), Upper Saddle River, New Jersey: Prentice Hall, ISBN 0-13-790395-2. Date Access: 06/23/2014.
- [7] GLOVER, F. e KOCHENBERGER, G. A. Handbook of Metaheuristics. Kluwer Academic Publishers, Boston. 2003.
- [8] Brownlee J. Clever Algorithms: Nature-Inspired Programming Recipes. 2012.
- [9] Dodge, Y. (2006) The Oxford Dictionary of Statistical Terms, OUP. ISBN 0-19-920613-9
- [10] Hornick M, Oracle Corporation, Java Specification Request 73, Java Data Mining (JDM), 2004.
- [11] <http://scikit-learn.org/stable/index.html>
- [12] <http://pymc-devs.github.io/pymc/>
- [13] E. Tapia, S. Intille, and K. Larson, Activity recognition in the home using simple and ubiquitous sensors, no. 1933. 2004, pp. 1–6.