**PAPER • OPEN ACCESS**

# Dimensionality reduction and clustering of time series for anomaly detection in a supermarket heating system

To cite this article: Lorenzo Salmina *et al* 2021 *J. Phys.: Conf. Ser.* **2042** 012027

View the article online for updates and enhancements.

# Dimensionality reduction and clustering of time series for anomaly detection in a supermarket heating system

**Lorenzo Salmina, Roberto Castello, Justine Stoll and Jean-Louis Scartezzini**

Solar Energy and Building Physics Laboratory, EPFL, Lausanne, Switzerland


Corresponding author: roberto.castello@epfl.ch

**Abstract**. A timely identification of an anomalous functioning of the energy system of an industrial building would increase the efficiency and the resilience of the energy infrastructure, beside reducing the economic wastage. This work has been inspired by the need of identifying, for a series of supermarket buildings in Switzerland, the failures happening in their heating systems across the years in an unsupervised and easy-to-visualize fashion for the building managers. The lack of any a-priori label differentiating between typical and anomalous behaviors calls for the usage of unsupervised machine learning methods to extract the relevant features to describe the system operations, to reduce the dimension of the feature space, and to cluster together similar patterns of operations. The method is validated on a standard supermarket building, where it successfully discriminates winter and summer operations from periods of refurbishment or malfunctioning of the heating system.

## 1. Introduction

The process of identifying which instances in a dataset stand out as being dissimilar to all others is called anomaly (or outlier) detection. Anomalies are common in many real-world applications, from identifying bank frauds, to medical records or fake news. This work has been inspired by the need of identifying, for a series of supermarket buildings in Switzerland, the irregularities happening in their heating systems across the years in an automated and easy-to-visualize fashion for building managers. The main difficulty in addressing such a problem is the lack of any a-priori labels differentiating between typical and anomalous behaviors of the system. Notoriously, unsupervised statistical learning tools, like clustering methods, are widely used to circumvent this problem. They are able to group and classify data points in a multidimensional feature space based on affinity-, proximity- or density-based metrics, and to isolate outliers which do not fit any of the clusters.

Being anomaly detection a broad topic, several attempts exist in the literature [1]. Frequently used anomaly detection methods for monitoring the operation of industrial components include parametric statistical modeling (autoregressive–moving-average or ARMA models), deep neural networks (like Convolutional Neural Networks and Auto Encoders) [2] and ruled-based algorithms. The choice of the detection method is highly dependent from the type of data used as input. Noisy or non-stationary dataset, as well as the availability of labelled data for training and validation of the algorithm, pose a major challenge. When it comes to the problem of classifying time series, unsupervised approaches like clustering are preferred as they do not require pre-defined patterns of anomaly to be defined neither the stationarity of the series [3]. Particularly interesting is the work of Ali et al. [4], where they present a workflow to reduce the dimensions of long and non-stationary multivariate time series and to cluster
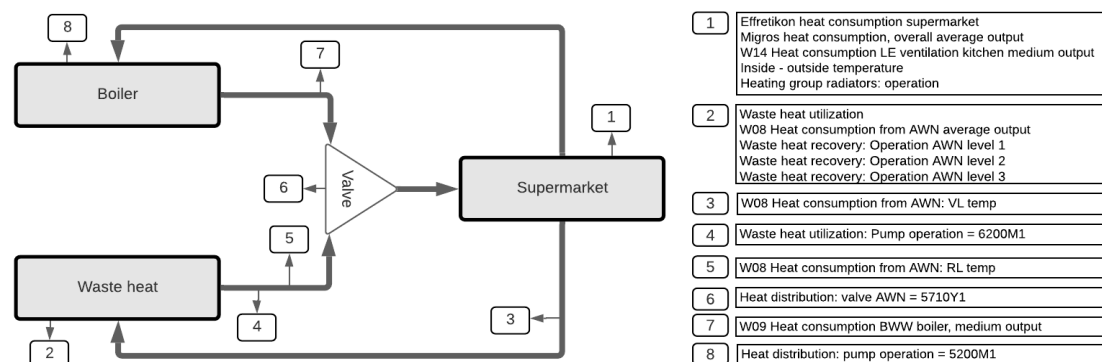
together recurrent temporal behaviors for pattern discovery and outlier detection. The case study they present on multivariate time series (recorded by a tri-axial accelerometer for monitoring the behavior of wild birds) well adapts to our case, where the outputs of multiple sensors are continuously recorded during the operation of the supermarket heating system.

In this work we aim at using dimensionality reduction techniques to extract the critical features needed to represent the operational status of the heating system along a week. We use these features to cluster together, in an unsupervised way, weekly patterns which are similar. Weeks which are not assigned to any class are potential candidates for non-standard operations of the system (or so-called anomalies) and should be hence scrutinized by building energy managers. The work is divided into the following parts. Firstly, we present a simplified, qualitative description of the heating system of the supermarket selected for the case study and the set of recorded sensor outputs. Secondly, we outline the methods for representing time series data and reducing their dimensions, through the combination of autoencoders and dimensionality reduction techniques. Thirdly, we describe the outcome of the clustering method and the classes found. We discuss the strength of the method by validating the results on the basis of the expert's feedback and we suggest future extensions of the method.

## 2. Data and Method

We analyze the heating system of a building owned by one of the Swiss supermarket chains located in the canton of Zurich. Based on one of the latest energy-saving strategies for heat recovery, the system of this supermarket combines the waste heat emitted by the condensers of the cooling systems (mostly refrigerators) with an electric boiler to heat the service water for routine operations. The system is controlled by a set of valves which enables or disables the contributions of each subsystem with the aim of keeping, through radiators, the indoor temperature at the desired level. The supermarket is equipped with several sensors which record information like indoor and outdoor temperatures, overall energy consumption, and state and percentage of activation of the valves at different nodes of the system. Figure 1 shows a schematic representation of the heating system.



| 1 | Effretikon heat consumption supermarket<br>Migros heat consumption, overall average output<br>W14 Heat consumption LE ventilation kitchen medium output<br>Inside - outside temperature<br>Heating group radiators: operation |
| 2 | Waste heat utilization<br>W08 Heat consumption from AWN average output<br>Waste heat recovery: Operation AWN level 1<br>Waste heat recovery: Operation AWN level 2<br>Waste heat recovery: Operation AWN level 3 |
| 3 | W08 Heat consumption from AWN: VL temp |
| 4 | Waste heat utilization: Pump operation = 6200M1 |
| 5 | W08 Heat consumption from AWN: RL temp |
| 6 | Heat distribution: valve AWN = 5710Y1 |
| 7 | W09 Heat consumption BWW boiler, medium output |
| 8 | Heat distribution: pump operation = 5200M1 |

**Figure 1.** Simplified outline of the supermarket heating system and the 16 variables measured by the sensors. Thick arrows represent the flow of water/air, while thin ones represent the location number of each sensor. Each location is mapped to the name and acronym of the measured variable.
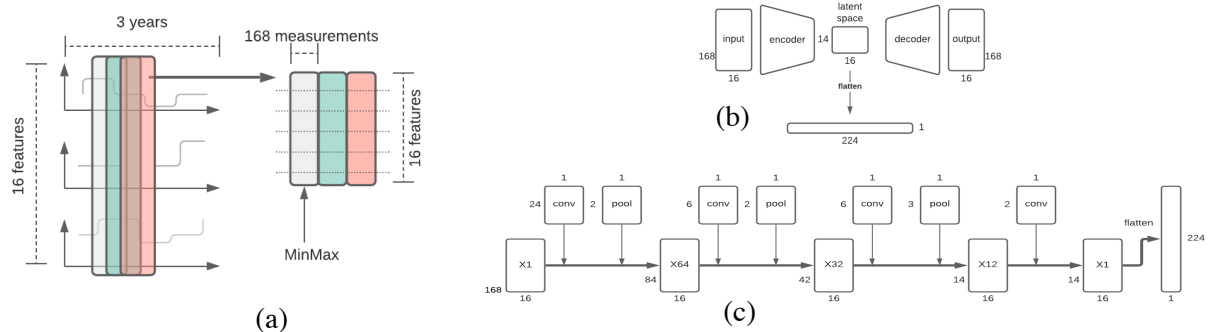
The standard operation mode of this heating system relies on waste heat from refrigerators. In case of shortage or extreme needs, the boiler is activated in order to keep the supermarket at the desired indoor temperature. As visible in Figure 1, there are several measurement points along the waste heat unit (the group "2" of the diagram). The variables are monitored for a period of almost four years, from June 2017 to May 2021 with a frequency of 15 minutes. We separate the data in two groups, before and after 31/05/2019. The first half (*training* set) is used to tune the model's workflow as it contains, according to the supermarket technicians, almost no anomalies of the heating system operations. The second half (*validation* set) is used to quantify the model's performance in detecting new anomalies in unseen data.

The method described in the following is applicable to any supermarket which operates through a similar heating system technology.

### 2.1. Feature selection

The first step consists of selecting a subset of features which is representative of the weekly operation status of the heating system described in Figure 1. Therefore, we compute the correlation matrix of all the variables recorded by the sensors, which reveals groups of highly positively as well as highly negatively correlated variables. After consultation with the supermarket building managers, we decide to keep only one feature for each highly correlated/uncorrelated group, and we include in the list mostly features closely related to the waste heat recovery system and valves. We end up with a selection of 16, out of the initial 30 features. We resample the time series to get a frequency of one measurement per hour: if the time series is continuous, we take the mean value over the hour, while if the time series is a sequence of binary signals (e.g., for the valve sensors the output is 1 if the valve is open and 0 if closed) we set the value to the number of times the valve is open within the hour.

Following the approach at [4] for the multivariate time series, we apply sliding windows of 168 hours length (one week) with stride of 1 hour on the selected set of 16 features over the two years covered by the *training* dataset. This results in a dataset of 17229 two-dimensional arrays, each one with 168 rows (one for each hourly-aggregated measurement) and 16 columns (one for each feature). The sliding window method is described in Figure 2a.



**Figure 2.** Schematic view of the sliding window mechanism for multivariate time series (a); Autoencoder network used to map features into the low-dimensional latent space representation (b); Detailed architecture of the encoder part of the network (c).

### 2.2. Dimensionality reduction

In this second step, the resulting tensor from the sliding window method is treated as a sample of points in a high-dimensional space. We use data reduction techniques to provide an alternative visualization and exploration of the time series behavior in a two-dimensional space. First, we trained a deep convolutional auto-encoder (DCAE) [5]. It belongs to a family of unsupervised models which maps the inputs into a low-dimensional representation (*latent space*) through fully connected layers embedded into an encoding and a decoding part, as shown in Figure 2b. The DCAE is trained in order to minimize the reconstruction error between the input and the output. To reach this, we split the original dataset of 17229 two-dimensional arrays in 90% for the training and the remaining 10% for validation. Each column of each 2D array is individually normalized using the min-max scaling. We use a training batch size of 100 arrays. We trained over 120 epochs using Adam optimizer, until the loss function (Mean Squared Error between the input and the output arrays) reaches a plateau. The final architecture of the encoder is visualized in Figure 2c.

We exploit the DCAE encoder to reduce the feature space from 16x168 to 16x14. It is worth to notice that the convolution and the max pooling are always performed column-wise: we opted for this approach in order to preserve the initial number of physical features (16) and to compress only the temporal ones

in the latent space representation. A post-training visualization of a typical input and the corresponding compressed representation and output is shown in Figure 3. The reconstructed latent space is a blurry version of the input image and it preserves the main patterns. Dark and bright regions are consistent in both input and output images. The low value of the MSE, obtained after the training over 120 epochs, together with the visual inspection of the decoded output images, make us confident that the latent space learned by the DCAE is a fairly good approximation of the input one.



**Figure 3.** Visualization of input, output and compressed representation from the DCA model for a typical week
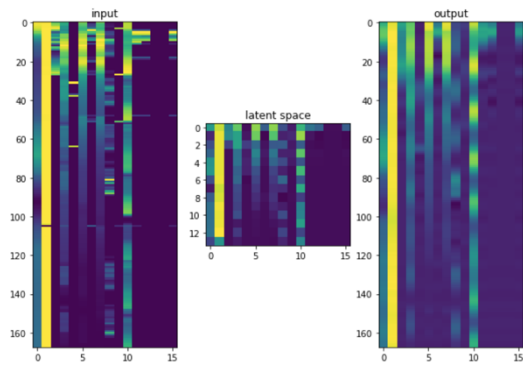
In order to further reduce the dimension of the latent space, the output after the encoder last layer is flattened and projected into a 2D space by means of linear and non-linear techniques, like Principal Component Analysis (PCA) [6] and Uniform Manifold Approximation and Projection (UMAP) [7]. We found that PCA provides low explained total variance (less than 50%) by using the first two components. We therefore explore the UMAP algorithm and we perform a grid search to tune its parameters, namely the number of components (*nr_comp*), the minimum distance (*min_dist*) and number of neighbors (*nr_neigh*). We choose the configuration (*nr_comp*=2, *min_dist*=0.4 and *nr_neigh*=60) which best separate in the two-dimensional plane winter-like from summer-like calendar weeks. In fact, during these two seasons the heating system behaviors are expected to be different and this should be reflected in the spatial proximity in the 2D plane of same-season weeks. This result is key in the tuning of the model parameters and it confirms how the chosen representation is sufficient to capture diametric status of the supermarket heating system.

*2.3.    Clustering method*

The third and final step consists into identifying clusters in the two-dimensional space defined in Section 2.2. A hierarchical density-based clustering method (HDBSCAN) [8] is used to identify the significant clusters in a fully unsupervised manner. It requires the tuning of one parameter only, which is the minimum size of the cluster itself. In general, the smaller this parameter is, the greater the number of clusters found. Conversely, the bigger the parameter, the more the clusters will be merged together, thus reducing the overall final number of clusters. We chose to set the minimum cluster size to 60 points. This choice ensures that the smallest cluster is driven by the trend of the supermarket heating systems over at least two and a half days (each point is a one-week window with 1-hour stride with respect to the closer point). In order to check the robustness of the method, we increase the minimum size without observing significant changes in the number of clusters found, pointing to the fact that the clusters are sufficiently dense and well separated in the chosen two-dimensional plane. For our dataset, HDBSCAN is preferred compared to other clustering methods such as K-means, which instead requires the number of clusters as a parameter, an information which is difficult to be known in advance (particularly for large datasets).
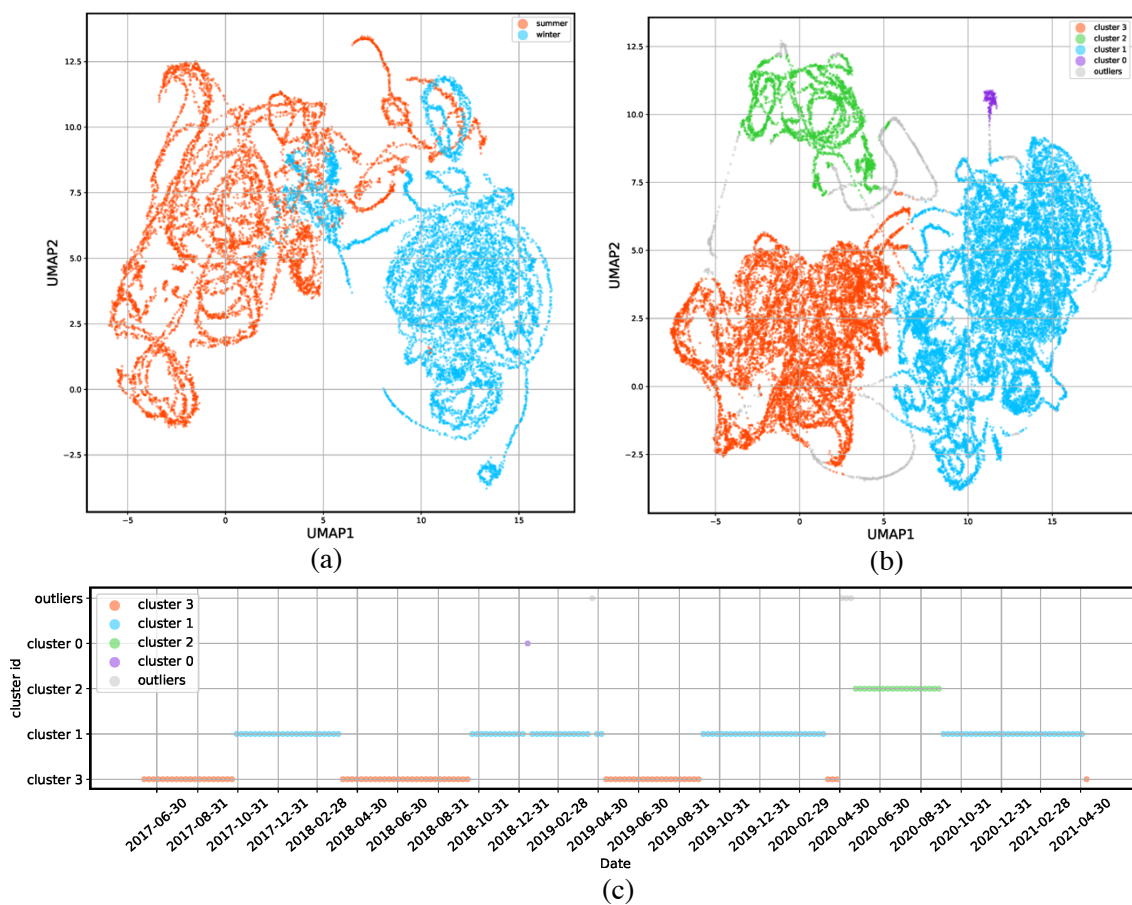
**3.  Results and discussion**

We project the evolution of the multivariate time-series of the supermarket heating system on a two-dimensional plane defined by the resulting features of the UMAP technique reported in Section 2.2. Figure 4a shows the evolution of the so-called *training* set. The DCAE and UMAP have been trained and tuned on this dataset in order to maximize the separation between the two larger heating periods, namely the summer-like and winter-like seasons of the year. As visible in Figure 4a, points (weeks) belonging to the same season are well separated in the 2D space defined by the feature extraction and dimensionality reduction, despite some spatial overlaps are present.

We test the model by injecting the *validation* set, namely a new set of time series recorded between May 2019 and May 2021. The sliding window technique described in Section 2.1 is applied to the *validation* set (16796 points) and to the *training* set again, and a low-dimensional representation is extracted through DCAE and UMAP for the *training* and *validation* set. The evolution of the combined dataset in the 2D space is visible in Figure 4b. The HDBSCAN clustering method is applied to bundle together similar weekly behaviors of the heating system with the scope of identifying anomalies. As a result of that, four classes of clusters and one class of outliers appear. We trace back the two largest clusters (cluster 1 and cluster 3) to weeks belonging to the summer-like and winter-like seasons. This confirms that the model is able to discriminate between seasonal behaviors of the supermarket heating systems. For the analysis of the other clusters, we proceed with a validation which uses a combination of calendar weeks and feedback from building managers.



**Figure 4.** Visualization of the weekly heating system operations for the *training* set (a) after UMAP. No clustering is applied and points are colored according to their season (summer-like in orange and winter-like in blue). Clusters found by HDBSCAN on the *training* and *validation* set (b). Clusters assignment to calendar weeks (c).

In Figure 4c we assign each calendar week of the four years covered by the dataset to one of the clusters found by HDBSCAN. Given that each week is represented by 168 sliding windows at the most and each window (corresponding to a point in the 2D plane) is assigned to a cluster, we decided to assign a cluster color to each calendar week according to the most recurrent cluster color among the sliding windows of that week. As a result of that, some weeks do not belong to the summer-like and winter-like clusters (*clusters 3* and *1*), pointing to different behaviors to be investigated. Particularly, the period from April to August 2020 could be ascribed to a systematic and prolonged operational mode of the supermarket heating system which induced atypical patterns in the monitored variables. Also, possible

effects due to the COVID-19 confinement on the supermarket energy regime cannot be excluded. To provide a qualitative validation of the model performance in classifying anomalous weeks, we ask the supermarket building managers to manually inspect weeks during the four years period, by indicating whether or not they recognize a faulty behavior of the system. As a result of this validation, we can ascribe the weeks from April to August 2020 (*cluster 2*) to a real period of refurbishment of the supermarket heating system. Moreover, the week belonging to *cluster 0* points to a real problem in the system, with a four times energy consumption increase with respect to the previous week. According to the logbook entries for that period, the problem has been fixed after one week, by changing the operative temperature of a thermovalve. No unusual or faulty behavior of the heating system has been found for the few weeks classified as outliers. They might be interpreted as false positive alarms, however, as visible in Figure 4b, they rather correspond to transition periods between seasons or to periods preceding anomalies, likely as an artifact of the sliding windows technique. As a consequence of that, we chose to consider all the weeks classified in this category as standard ones.

## 4. Conclusions

In this work we present a machine-learning based method to identify, for a supermarket in Switzerland, potential failures happening to its heating systems across the years in an unsupervised and easy-to-visualize fashion for the building managers. The lack of any a-priori labels differentiating between typical and anomalous behaviors calls for the usage of unsupervised methods to extract the relevant features to describe the system operations, to reduce the dimension of the feature space, and to cluster together similar patterns of operations. The method successfully discriminates winter from summer operations and potentially also periods of refurbishment and malfunctioning of the heating system. The method is scalable to all those supermarkets where the heating system operations can be described by the reduced set of features used in this study and it has the advantage of not requiring any manual a-priori labelling of the anomalous behaviors of the system.

## References

[1]    Chandola, V.; Banerjee, A. & Kumar, V., Anomaly detection: A survey, ACM Computing Surveys (CSUR) 41 (3), 15 (2009), https://doi.org/10.1145/1541880.1541882

[2]    G. Pang, C. Shen, L. Cao, and A. Van Den Hengel, Deep Learning for Anomaly Detection: A Review. ACM Comput. Surv. 54, 2, Article 38 (2021), https://doi.org/10.1145/3439950

[3]    S. Aghabozorgi, A. Seyed Shirkhorshidi, T. Ying Wah, Time-series clustering - A decade review, Information Systems, Volume 53, 2015, Pages 16-38, https://doi.org/10.1016/j.is.2015.04.007

[4]    Ali, M., Jones, M.W., Xie, X. et al. TimeCluster: dimension reduction applied to temporal data for visual analytics. Vis Comput 35, 1013–1026 (2019). https://doi.org/10.1007/s00371-019-01673-y

[5]    G. E. Hinton, R. R. Salakhutdinov, Reducing the Dimensionality of Data with Neural Networks, Science (2006) Vol. 313, Issue 5786, pp. 504-507 https://doi.org/10.1126/science.1127647

[6]    Jolliffe, I. T., and Cadima, J. Principal component analysis: a review and recent developments, Philosophical transactions, (2016) Series A, Mathematical, physical, and engineering sciences, 374(2065), 20150202. https://doi.org/10.1098/rsta.2015.0202

[7]    McInnes L., Healy J. and Melville J., UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, 2020, https://arxiv.org/abs/1802.03426v3

[8]    Campello R.J.G.B., Moulavi D., Sander J. (2013) Density-Based Clustering Based on Hierarchical Density Estimates. In: Pei J., Tseng V.S., Cao L., Motoda H., Xu G. (eds) Advances in Knowledge Discovery and Data Mining. PAKDD 2013. Lecture Notes in Computer Science, vol 7819. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-37456-2_14