

Clustering Activities of Daily Living without Supervision using DBSCAN

Jeremy Leyden

Computer Science Department – Machine Learning and Deep Learning

Kent State University

Kent, Ohio

jermo@leydenfamily.com

Abstract—This study explores the application of unsupervised clustering techniques to smart home sensor data for identifying some sample Activities of Daily Living (ADLs) without labeled supervision. Using the publicly available CASAS Ordonez A and B datasets, which record timestamped activations of various ambient sensors in a two-person household, a set of contextual and time-based features was engineered at the event level. These features include sensor type, location, activation frequency, time-of-day encoding, inter-event timing, and room-level mappings, all aimed at capturing behavioral patterns indicative of daily routines. The overall clustering performance was evaluated using external metrics—Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI)—against ground-truth activity labels. Among the clustering algorithms tested, DBSCAN achieved the most promising results (ARI = 0.320, NMI = 0.399), outperforming K-Means and other baseline approaches, likely due to its robustness to noise and ability to model complex clusters that don’t solely rely on Euclidean distance. These preliminary findings suggest that unsupervised approaches, when combined with thoughtful feature engineering, can meaningfully uncover latent structure in smart home sensor data, offering potential for real-world activity recognition without manual annotation.

Keywords—ADLs, supervised learning, feature engineering, sensors, ARI, NMI, clustering, DBSCAN

I. INTRODUCTION

As populations age and the demand for independent living increases, intelligent systems for monitoring Activities of Daily Living (ADLs) in smart home environments have become increasingly normalized. ADLs—such as eating, sleeping, grooming, and toileting—serve as key indicators of health, lifestyle, and cognitive function. While supervised machine learning has been widely used to detect ADLs from sensor data, these approaches typically rely on time-consuming and labor-intensive manual labeling, which limits scalability and real-world deployment. Unsupervised learning presents a compelling alternative, offering the potential to uncover meaningful patterns in sensor data without requiring ground-truth annotations. However, clustering in such settings is inherently challenging due to the complexity and variability of human behavior, sparsity and noise in sensor activations, and the difficulty of defining appropriate feature representations.

This work aims to investigate the feasibility of using unsupervised clustering to recognize ADLs from ambient sensor data collected in a residential setting. By utilizing the Ordonez A and B datasets, which contain rich streams of timestamped

sensor activations across multiple rooms and sensor types, cluster predictions can be performed from temporal features. To capture these temporal and spatial dynamics of human activity, a set of features is engineered at the event level, including activation timing, sensor transitions, room-level occupancy patterns, and time-of-day encodings. Following that, this study applies density-based and centroid-based clustering methods, namely DBSCAN, and evaluate its performance using external cluster validity indices, specifically Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI), based on the available labeled data. Our results demonstrate that DBSCAN, with its ability to handle noise and detect clusters of arbitrary shape, achieves stronger performance, highlighting the importance of algorithm choice in this domain.

This study contributes to the growing field of unsupervised activity recognition by demonstrating how careful feature design and appropriate clustering methods can reveal latent structure in smart home sensor data, offering a path toward scalable and unobtrusive ADL monitoring.

II. METHODOLOGY

A. Datasets and Preprocessing

This section describes the complete data preprocessing pipeline, feature engineering strategy, and clustering methods used for unsupervised ADL recognition using the Ordonez A and B smart home datasets. Our approach focuses on modeling the sensor activation stream as a sequence of related events, from which we derive a set of temporal, spatial, and contextual features to support unsupervised learned and clustering. The experiments utilize two datasets acquired from Canvas, `OrdonezA_Sensors.txt` and `OrdonezB_Sensors.txt`, each corresponding to a different individual living independently in a multi-room smart home environment. Both datasets record ambient sensor activations over several weeks, covering a variety of common activities of daily living. Each event in the dataset includes start and end timestamps, sensor location, type, and place (e.g., kitchen, bathroom). In addition to sensor data, we also incorporate manually labeled ADLs provided in `OrdonezA_ADLS.txt` and `OrdonezB_ADLS.txt`, which specify the start and end times for each activity. These labels are not used during clustering but serve to evaluate clustering performance afterwards.

B. Temporal Relationships

To work with temporal dynamics, we first parsed the start and end timestamps of each sensor activation into datetime objects and computed each event's duration in seconds. Events with missing timestamps were discarded. This step ensures that each data point represents a valid time-bound sensor interaction. The datasets from both users were merged into a unified event stream. To retain the identity of each subject, a Set column was added (denoting 'A' or 'B'). The merged data was sorted chronologically by activation time. Additionally, we computed the time since the last activation, which captures temporal gaps between successive events, and filled any missing values with zero for initialization.

C. Spatial/Transition Based Relationships

Each sensor's location (e.g., "Fridge", "Seat") was mapped to a higher-level room category (e.g., "Kitchen", "Living Room") using a predefined dictionary. The resulting Room feature captures spatial context and was encoded numerically using category encoding. This abstraction facilitates understanding behavioral patterns at the room level rather than individual sensor granularity. To capture temporal dependencies and activity transitions, we derived several features based on sensor sequences, such as the type of the previous sensor event in the stream (PrevSensorType), the categorical combination of previous and current sensor types (SensorTransition), and a binary flag indicating whether the current sensor type differs from the previous one (SensorTypeChange). Such transition-based features are essential in modeling the flow of activities and recognizing common routines.

D. Other Features

Two features were engineered to quantify behavioral patterns over time, which are the average time interval between consecutive activations of the same sensor location. This reflects how frequently a space is used (SensorFrequency), and the cumulative time gaps between events in the same room, serving as a proxy for time spent in each space (RoomOccupancy). These features aim to capture user routines and dwell time in functional areas of the house.

To further enhance the feature space, one-hot encoding was applied to the 'Type' of sensor, introducing a binary indicator for each unique sensor type (e.g., PIR, Magnetic, Electric, Pressure). This allowed the model to treat sensor types independently rather than ordinally. In addition, a SensorActivationCount feature tracks the number of times each sensor location has been triggered up to a given point. This helps identify areas of frequent usage and may correlate with specific activities like cooking or grooming.

E. Clustering and Evaluation Metrics

While the clustering model operates in an unsupervised manner, the ground-truth ADL labels were assigned post hoc by matching each sensor activation to its corresponding activity window from the labeled ADL files. This labeling was done per timestamp and per subject, enabling evaluation of clustering quality using standard metrics. All numerical features were scaled using a StandardScaler, which centers each feature to zero mean and unit variance. This step is critical to prevent

features with large numeric ranges (e.g., durations, time gaps) from dominating the clustering algorithm. From this, the main clustering algorithm evaluated here is DBSCAN (Density-Based Spatial Clustering of Applications with Noise): A density-based method that discovers clusters of arbitrary shape and can handle outliers (noise). Hyperparameters `eps` and `min_samples` were tuned based on ARI/NMI performance. The following external clustering evaluation metrics were used: Adjusted Rand Index (ARI), which measures the similarity between the clustering result and the ground truth, adjusted for chance, and Normalized Mutual Information (NMI), which quantifies the mutual dependence between cluster assignments and true labels, normalized between 0 and 1. These metrics allow assessment on how well the clustering aligns with the true activity structure, even though labels were not used during clustering.

III. RESULTS

DBSCAN was used for clustering at an `eps` value of 2.2, with a respective `num_samples` value of 4. While many values were tested using grid-search and checking manually within a for-loop, these values were deemed optimal for the feature dataframe provided. Once clustering was performed, the best ARI/NMI values were achieved at 0.320 and 0.399 respectively. This data is shown in Table 1.

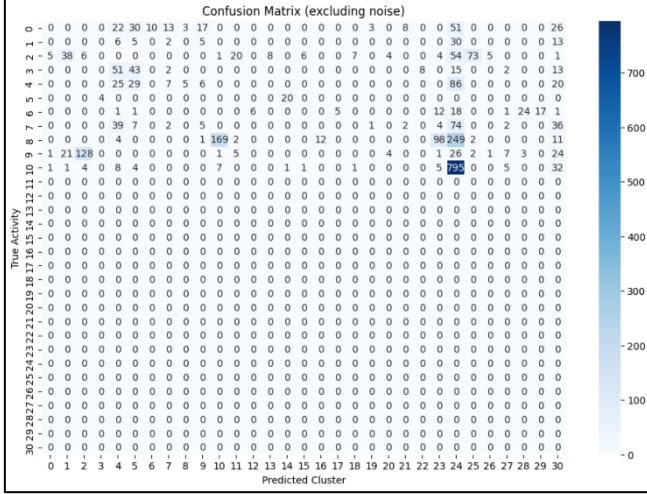
Table 1: Best Results

ARI:	0.320
NMI:	0.399
Eps:	2.2
Num_samples:	4
Clustering Method:	DBSCAN

This table displays the results of clustering with DBSCAN at the correspond `eps` and `num_samples` values. These values were chosen after examining all possible `eps/num_sample` combinations between 1 and 40.

In order to display the results of which predictions were correct/incorrect in regards to the clustering attempt, a confusion matrix was generated, displaying the true activity plotted against the predicted cluster, with an intensity showing color for its frequency/value to represent the number displayed. This plot is displayed in Figure 1.

Figure 1: Confusion Matrix



This figure represents the accuracy of the predicted clusters against the true labels. This plot uses DBSCAN and excludes noise points, labeled as -1.

Following the generation of the confusion matrix, PCA was performed to better understand the relationship between each feature and the given results. As a result, the use of PCA flattens the dimensions into PCA1 and PCA2, allowing for a 2D representation of the features observed. This figure is shown in Figure 2A, and can be compared to the ground-truth PCA shown in Figure 2B.

Figure 2A: PCA for Predicted Clusters



This figure represents the feature spread of the predicted clusters, for comparison with the ground-truth PCA,

Figure 2A: PCA for Ground-Truth Labels



This figure represents the feature spread of the ground-truth labels, for comparison with the predicted cluster PCA,

IV. CONCLUSIONS

This study explored the effectiveness of unsupervised clustering for inferring Activities of Daily Living (ADLs) from ambient sensor data using the Ordenez A and B datasets. The core objective was to determine whether meaningful patterns of human behavior could be identified directly from raw sensor activations, without relying on supervised labels during training. Through extensive preprocessing and feature engineering—incorporating temporal, spatial, and contextual elements such as sensor type transitions, time since last activation, and cyclic time encodings—we constructed a rich feature space capturing the underlying structure of daily routines. The density-based DBSCAN algorithm was chosen for its ability to detect clusters of arbitrary shape and to ignore outliers, making it well suited to the noisy and event-driven nature of smart home data. This experiment revealed that DBSCAN was capable of clustering sensor events with moderate alignment to ground-truth ADLs, achieving an Adjusted Rand Index (ARI) of 0.320 and a Normalized Mutual Information (NMI) of 0.399 under optimal parameters ($\epsilon = 2.2$, $\text{min_samples} = 4$). These metrics reflect the partial discovery of latent activity patterns purely from sensor dynamics. While unsupervised performance naturally lags behind supervised methods, the results demonstrate that structure does exist within the sensor data stream that aligns with real-world behavior. The interpretability of clusters and their relationship to known ADLs also highlights the value of well-designed features and context-aware preprocessing. Future work will focus on enhancing the representation of transitions between activities, incorporating temporal segmentation windows, and exploring hybrid models that bridge unsupervised embeddings with weak supervision or semi-supervised refinement. Additionally, further evaluation on alternative smart home datasets will help assess the generalizability of the DBSCAN-based approach.

V. REFERENCES

- [1] Rodriguez, Mayra Z et al. "Clustering algorithms: A comparative approach." PloS one vol. 14,1 e0210236. 15 Jan. 2019, doi:10.1371/journal.pone.0210236.
- [2] Ming-Syan Chen, Jiawei Han and P. S. Yu, "Data mining: an overview from a database perspective," in IEEE Transactions on Knowledge and Data Engineering, vol. 8, no. 6, pp. 866-883, Dec. 1996, doi: 10.1109/69.553155.
- [3] Rangaprakash, D et al. "Density-based clustering of static and dynamic functional MRI connectivity features obtained from subjects with cognitive impairment." Brain informatics vol. 7,1 19. 26 Nov. 2020, doi:10.1186/s40708-020-00120-2.
- [4] D. Deng, "DBSCAN Clustering Algorithm Based on Density," 2020 7th International Forum on Electrical Engineering and Automation (IFEEA), Hefei, China, 2020, pp. 949-953, doi: 10.1109/IFEEA51475.2020.00199.