Using Continuous Sensor Data to Formalize a Model of In-Home Activity Patterns

Abstract. Formal modeling and analysis of human behavior can properly advance disciplines ranging from 3 psychology to economics. The ability to perform such modeling has been limited by a lack of ecologically-4 valid data collected regarding human daily activity. We propose a formal model of indoor routine behavior 5 based on data from automatically-sensed and recognized activities. A mechanistic description of behavior 6 7 patterns for identical activity is offered to both investigate behavioral norms with 99 smart homes and 8 compare these norms between subgroups. We identify and model the patterns of human behaviors based on 9 inter-arrival times, the time interval between two successive activities, for selected activity classes in the smart home dataset with diverse participants. We also explore the inter-arrival times of sequence of activities 10 11 in one smart home. To demonstrate the impact such analysis can have on other disciplines, we use this same smart home data to examine the relationship between the formal model and resident health status. Our study 12 reveals that human indoor activities can be described by non-Poisson processes and that the corresponding 13 distribution of activity inter-arrival times follows a Pareto distribution. We further discover that the 14 combination of activities in certain subgroups can be described by multivariate Pareto distributions. These 15 findings will help researchers understand indoor activity routine patterns and develop more sophisticated 16 models of predicting routine behaviors and their timings. Eventually, the findings may also be used to 17 automate diagnoses and design customized behavioral interventions by providing activity-anticipatory 18 services that will benefit both caregivers and patients. 19

20 KEYWORDS

Human dynamics, Population modeling, Pareto distribution, Pervasive environment, Activity recognition

22 1 INTRODUCTION

The wealth of data that can now be collected by 23 ambient sensors facilitates the development of new 24 models of human behavior supported by empirical 25 evidence. In this paper, we propose formal models of 26 human activities for indoor environments. Specifically, 27 we analyze and model the sequences and timings of basic 28 everyday activities for smart home residents. Offering 29 such models provides a basis for making claims 30 31 regarding human behavior and differentiating behavior 32 strategies for population subgroups (healthy, dementia). We validate our models by using multiple years of 33 ambient sensor data collected in smart homes. We find 34 that activity arrival rates can be mathematically modeled 35 and that model parameters differ between healthy older 36 adults and older adults with chronic health issues. These 37 analyses allow researchers to better understand the 38 impact of health conditions on routine behavior and can 39 be used to predict diagnosis categories for individuals 40 41 based on automatically-sensed activity patterns. 42 Due to limitations with real-world data collection methods, previous models for human activities did not 43 provide sufficient information about the dynamic 44 property of human behaviors. They typically assumed 45 that human activities can be modeled by Poisson 46 processes and that the inter-arrival time, or the time 47 interval between two successive activities, follows an 48

exponential distribution. This assumption models 4۵ activities as occurring at a constant rate [1]-[4]. 50 51 However, this model does not capture the fluctuation 52 that may occur in activity arrival rates. With the advent of quantifiable mobile data that can be collected 53 unobtrusively and continuously, researchers recently 54 proposed the use of heavy-tailed distributions to 55 56 describe human dynamics [5]–[9].

Our approach in this paper is to build a general model 57 58 of human activities that involves real-time data 59 collection in everyday environments based on ambient sensor data collected in smart homes. We perform our 60 61 data collection on subjects inside their own homes. The data collection reflects routine human behavior without 62 requiring any alteration to the environment or activities, 63 facilitating an ecologically valid analysis. We analyze the 64 inter-arrival times of automatically-labelled smart home 65 sensor data (e.g., cooking, eating), and find activity 66 interdependencies in subgroups (healthy older adults 67 and older adults with chronic health problems). To 68 investigate the relationship between behavior changes 69 and health problems, we use a case study with 65 months 70 71 of data from one smart home. This behavior-driven 72 sensor data shows that activity routines can be modeled by non-Poisson processes. The activity inter-arrival 73 times follow a heavy-tailed distribution, specifically a 74 Pareto distribution. 75

We find that model parameters for activity arrivalrates differ between healthy older adults and older

adults with health issues. The resulting mathematical 1 models open up the possibility of recognizing the 2 development of health problems and providing efficient 3 interventions and assistance. Once differences in 4 patterns among subgroups are found, they can be used 5 6 to better understand the impact of culture, age, and education on daily routines. The design of technology-7 based tools such as agent- and human-oriented software 8 9 and hardware systems [8], [10]-[12] can also greatly benefit from this work. Researchers in the fields of 10 sociology, psychology, and anthropology will also be able 11 to align their studies with customized and personalized 12 healthcare systems. 13

Our study provides evidence to support three 14 hypotheses of human routine behaviors in home 15 16 environments. First, human behavior can be described 17 by formal statistical distributions. Second, data 18 supporting this conclusion can be collected using 19 ambient sensors in an ecologically valid manner. Third, the Pareto model and its properties, such as the 80/20 20 rule, can be useful for the study of human dynamics and 21 investigation of hypotheses because of its ability to 22 model human behavior patterns. Our study first analyzes 23 and models inter-arrival times of identical indoor 24 behaviors based on both 99 smart homes and subgroups 25 26 of older adults. We further study the inter-arrival times 27 of activity sequences from one smart home. The findings 28 of this study will offer the potential to automate 29 diagnoses and design customized behavioral 30 interventions.

31 2. LITERATURE REVIEW

32 Maturing pervasive computing technologies have sparked a new wave of human behavior analysis and 33 34 resulted in new theories regarding human behavior 35 patterns. Barabási's study of the timing of consecutive electronic and physical mail messages sparked a model 36 37 of human dynamics as a heavy-tailed distribution [5], [13]. A queuing model and heavy-tailed distribution 38 were introduced in Barabási's study to explain the large 39 time gap between sent messages after a burst of 40 responses. 41

After Barabási's discovery, scientists use heavy-42 tailed distributions to explain human behavior in diverse 43 44 domains, ranging from social science to health care. In 45 the social network field, heavy-tailed distributions are used to characterize the dynamics of popularity based on 46 47 diverse digital platforms, such as Wikipedia, blog posts, 48 Android applications, Web pages, and Twitter [14]–[20]. As an example, Li et al. show that the behavior-based 49 popularity of Android applications follows the Pareto 50 principle [17]. Tsompanidis et al. also discover that web 51 traffic flow size can be explained by the Pareto 52 distribution [19]. Similarly, researchers presented a list 53

of social and organizational power laws, one kind of
heavy-tailed distribution, to describe human behavior
[21]-[23]. Specifically, the power law distribution
identifies the number of inter-firm relationships
observed from linkages between firms: suppliers,
customers, and owners [22], [23].

60 Further, scientists use heavy-tailed distributions to model and predict human mobility [24]-[30]. For 61 62 example, GPS-based human movement patterns can be captured by heavy-tailed flights for different 63 transportation modes, including walking/running and 64 car/taxi [28]-[30]. Regardless of transportation modes, 65 the distribution of user's moving distances, from visited 66 locations to the target location, can be modelled by the 67 Pareto distribution [27]. 68

69 Besides heavy-tailed distributions, other 70 mathematical models are also used to uderstand a more 71 varied set of human activities than basic movements. A 72 mixture of Gaussian intensities model is introduced to explain activities, such as exercising and eating, that have 73 time-varying, interdependent, and periodic properties 74 [31]. The temporal granularity algorithm, considering 75 behaviors happened within a time interval instead of at 76 an exact timestamp, is used to identify frequent 77 behavioral patterns, such as receiving a call, 78 79 sending/receiving a text message, and holding a meeting 80 [32].

81 In addition to mathematical formalisms, researchers 82 adopt machine learning methods to model aspects of 83 human behavior [33], [34]. For example, inverse 84 reinforcement learning (a method which flips the 85 problem of traditional reinforcement learning and learns an agent's rewards by observing its behavior) models 86 human driving routines to help aggressive drivers 87 improve their driving style [33] and to find taxi driver's 88 89 preferences on working regions and times [34].

90 Given the development of these diverse models to understand human behavior patterns, we propose to 91 92 extend previous work further by modeling indoor 93 behavior patterns based on ambient sensor data. 94 Although researchers have analyzed raw sensor data and 95 design features from the raw sensor data to understand human behaviors in smart environments [35]–[38], the 96 substantial number and diversity of raw sensor event 97 patterns are difficult to provide a rich vocabulary to 98 99 express human behavior. Analyzing the labeled activities 100 from sensor event sequences resolves the concern.

101 In this paper, we model generalized human behavior 102 based on ambient sensor data collected in smart homes. 103 Other work has similarly focused on labeling and 104 analyzing smart home-based behavior. Some of this prior 105 research introduces data-driven techniques for 106 recognizing or predicting daily activities in smart home 107 environments based on continuous sensor data [39]-[46]. Based on raw or activity-labeled sensor data, other 108

studies analyze and assess and individual's physical and 1 mental health stated associated with the observed 2 behavior [47], [48]. Outside of the health domain, 3 researchers have also analyzed behavior patterns from 4 ambient sensor data to predict the associated energy 5 6 consumption [49], a useful step in designing energyefficient automated buildings. However, these 7 techniques do not offer a mechanistic description of 8 9 indoor behavioral patterns. Furthermore, they have not yet attempted to describe behavior patterns at a 10 population level. 11

Our work explores several research problems. First, 12 we utilize activity recognition methods to label sensor 13 data in real time with corresponding activity labels. 14 Second, based on this labeled data, we analyze activity 15 16 inter-arrival times and construct heavy-tailed 17 distributions, specifically Pareto distributions, to 18 describe routine patterns for smart home residents in 19 everyday environments. Third, we investigate the patterns of selected activities both at a group level with 20 99 smart homes and between subgroups of older adults 21 (healthy, chronic condition) to have a generalized 22 understanding of behavior patterns and their differences 23 across a population. Fourth, we analyze the information 24 from our model to determine its value as an indication of 25 26 a person's health status.

27 3. METHODOLOGY

We propose a method to formally model activity 28 timings from behavior-based sensor data. First, we 29 30 monitor sensor events and automatically label the events with corresponding activity names using machine 31 learning models [41]. We then use change point 32 detection (CPD) [50] to segment data into sequences that 33 34 represent single, uninterrupted activities. Once the data 35 is segmented, we apply a well-known statistical method, extreme value theory (EVT) [28] to remove noise. We use 36

the remaining data to perform distribution fitting of the

histogram of activity inter-arrival times. We model the 53 data for 82 different probability distribution functions 54 (pdf) and determine the best distributions based on 55 minimizing the summation of the squared errors (SSE). 56 We utilize non-Poisson processes to model the inter-57 58 arrival times of human behavior routines and postulate that activity inter-arrival times can be approximated by 59 Pareto distributions. The modeling steps are illustrated 60 61 in Fig. 1.

We repeat these modeling steps for each activity 62 63 separately both for a complete sample of 99 smart home residents as well as for two subgroups of older adults 64 within the sample: healthy and chronic cognitive or 65 physical health conditions. We test the hypothesis that 66 differences in health status between subgroups may be 67 68 significantly reflected by patterns of each activity. 69 Analyzing inter-arrival times for one activity at a time is 70 reflective of analysis techniques used in previous studies 71 [51]–[54]. Furthermore, focusing on a single activity allows us to understand the potential relationship 72 73 between the activity of interest and population subgroups as well as identify differences in model 74 parameters between activities. On the other hand, 75 analyzing individual activity may prevent us from 76 holistically examining a person's entire behavior routine. 77 78 Thus, we additionally analyze the entire activity 79 sequence patterns for one of the smart homes.

80 3.1. Data Collection and Processing

In this study, we collect data from 99 smart homes to
investigate routine behavior patterns for selected
activities. We provide details on the first four steps of the
process in Fig. 1: data collection, activity labeling,
segmentation, and noise filtering.

86 3.1.1. Data Collection

Box (SHiB) [41], [55]. In each smart home, four types of
ambient sensors are installed: infrared motion, magnetic



52

37

(door) contact, light level, and temperature level. These 1 sensors are discrete event sensors and thus only 2 generate a message (sensor event) when there is a 3 4 change in a state, such as a refrigerator door opening or closing. The sensors are installed throughout the house 5 6 in each room including the kitchen, living room, dining room, bedrooms, bathrooms, office, and laundry room. 7 Infrared motion sensors include narrow-area and wide-8 9 area sensors. Narrow-area motion sensors detect heatbased movements within a one-meter diameter area. 10 They are attached to the ceiling above specific objects or 11 areas in the home, such as above a participant's favorite 12 chair or above a sink. Wide-area motion sensors perceive 13 movements occurring anywhere in an entire room. 14 These sensors are placed on ceiling corners in large 15 16 rooms, such as the dining or living room. Magnetic 17 sensors detect the use of doors and cabinets, such as in entering or leaving the home or accessing items within 18 19 kitchen or bathroom cabinets. Temperature sensors can be useful in detecting activities that change the heat level 20 in an area of the home, such as showering or cooking. 21 Similarly, light sensors can help us identify activities 22 occurring within a home as well as seasonal effects of 23 24 light levels. 3.1.2. Activity Labeling of Sensor Data 25 26 Activity labeling provides us with a rich vocabulary 27 to express human behavior. We employ automated 28 activity recognition techniques to label collected data 29 with eleven activity classes. The set of activities that we categorize and use in this analysis are seven activities: 30 Relax, Cook, Eat, Personal Hygiene, Wash Dishes, Sleep, 31

and Work. We use a separate class, Other Activity, to
 recognize unidentified sensor events.

We apply automated activity recognition (AR), a heavily-investigated challenge [31], [34]–[38], to map a sequence of captured sensor events onto one of the activity classes. Our AR steps are based on an approach that has been previously-validated for real-time activity recognition from streaming sensor data [39]. First, we

54

extract features from the raw data collected from the 55 discrete event sensors (see Fig. 2). We move a fixed-size 56 sliding window over the time-ordered sensor data and 57 compute feature values for each window [33], [39], [40]. 58 Within each window, sensor events may be widely 59 60 spread apart in time. To take this into account, a timebased weighting factor is applied to account for the 61 relative temporal distance between sensor events. 62 63 Second, after training a random forest classifier with ground truth pre-labeled sample data, the resulting 64 model can provide activity labels for data based on 65 features extracted from sensor sequences. The sequence 66 of sensor events in a window provides a context for 67 labeling the last (most recent) sensor event. The 68 approach we utilize is distinctive in that it is designed to 69 70 provide activity labels in real-time from ambient sensor 71 data collected continuously in real homes and to build a 72 generalizable model based on training ground truth-73 labeled data. After training, the resulting model can provide activity labels for data obtained in new smart 74 home settings. Our approach to activity recognition 75 76 yields an average of 95% accuracy and 0.78 F1 scores for activity labeling based on the three-fold cross-validation 77 78 method [41].

After applying AR, all unidentified sensor events are 79 80 assigned to the class, Other Activity. The drawback of 81 putting all undefined activities into a single Other 82 Activity class is that we cannot distinguish those 83 activities from each other. Each of them may shed light on the behavioral routines for the residents but because 84 they are grouped together, they cannot be analyzed as 85 individual activities. Therefore, for the Other Activity 86 category, we employ a k-means clustering algorithm to 87 divide the category into a specified number, k, of clusters. 88 The value of k is chosen using an elbow curve method, 89 which shows the minimized sum of squared distances of 90 samples from the closest cluster center. Thus, for each 91 92 home, our activity classes include both the seven 93 predefined and the cluster-generated activity classes, rk.

40				94	which are labeled cluster1 through cluste		
41							
42	A fixed	d-size sliding wi	ndow				
43					Features in each window		
44	2017-02-22 2017-02-22	11:42:50.57	FrontDoor Entry		Hour of day for current event Seconds since the start of the day for the current event		
45	2017-02-22	2017-02-22 11:42:51.57 Hall		Window duration			
46	2017-02-22 2017-02-22	11:42:51.99 11:42:53.01	Entry	try Seconds since try Most frequent mtDoor Feature Most frequent Current event	Seconds since previous event Most frequent event sensor in the previous window		
47	2017-02-22	11:42:54.82	FrontDoor		Most frequent event sensor in the window before that Current event sensor		
48	2017-02-22 2017-02-22	11:42:55.74	Entry Entry	Extraction	Most recent location sensor		
49	2017-02-22	11:42:57.18	Entry		Number of events in the window for each sensor Time since last fired for each sensor		
50	2017-02-22 2017-02-22	11:42:58.15	Hall		Maximum value of se Complexity of event	Maximum value of sensor Complexity of events in window	
51	2017-02-22	11:43:00.22	Kitchen		Number of motion events in window		
52	2017-02-22 2017-02-22	11:43:08.21 11:43:09.47	Kitchen Kitchen				
53							



4

The activity recognition algorithm labels each sensor 1 event with a corresponding activity class. The algorithm 2 does not, however, indicate the beginning or ending of 3 each activity occurrence. This information is valuable for 4 our analysis because we want to consider the inter-5 6 arrival times of each activity's start as part of a person's overall routine. To segment labeled data into individual 7 activities, we utilize an unsupervised method referred to 8 9 as Change Point Detection (CPD). CPD identifies the point in time where the state of the underlying process 10 changes [56], [57]. CPD can be used to detect real-time 11 activity transitions or changes in the data between two 12 successive windows of sensor events [58]. An example of 13 sensor events with corresponding times, activity labels, 14 and detected change points is shown in Fig. 3, where 15 16 transitions are indicated by horizontal lines. The first 17 event in an activity segment (the first sensor event after 18 a change point) is considered the start of an activity and 19 the last event in a segment (the last sensor event before a change point) represents the end of an activity. In this 20 study, we use a CPD method that is based on Bayesian 21 online learning [59]. Given the segmented data, we label 22 each segment with the most majority activity label for 23 that segment. The labeled activity in each segment 24 provides us the time and activity information and allows 25 26 us to calculate the inter-arrival times of each activity (in 27 hours), defined as the time between two successive start 28 times of the activity.

29	Change Point	Date and Time	Labelled Activity
20	0	2012-08-24 16:06:06	Other Activity
50	0	2012-08-24 16:06:07	Other Activity
31	1	2012-08-24 16:12:24	Other Activity
32	0	2012-08-24 16:12:26	Relax
52	0	2012-08-24 16:13:34	Relax
33	0	2012-08-24 16:13:35	Relax
3/	0	2012-08-24 16:34:52	Relax
54	0	2012-08-24 16:34:53	Relax
35	0	2012-08-24 16:35:05	Relax
26	0	2012-08-24 16:35:06	Relax
30	0	2012-08-24 16:35:06	Relax
37	0	2012-08-24 16:35:07	Relax
20	0	2012-08-24 16:35:10	Relax
38	0	2012-08-24 16:35:10	Relax
39	0	2012-08-24 16:35:13	Relax
40	0	2012-08-24 16:35:14	Relax
40	1	2012-08-24 16:35:17	Relax
41	0	2012-08-24 16:35:17	Relax
42			

Figure 3. A sample of CPD application to the sensor data.
A change point value of 1 indicates a transition/activity
start time, 0 indicates no change. The 0 right before the
next transition is the end time of an activity. A transition
is detected and shown by a horizontal line.

48

49 3.1.3. Participant Information

In addition to collecting sensor data for 99 smart homes, 50 we also store four additional parameters for each home: 51 the number of residents as well as resident ages, 52 53 education levels, and physical and mental health statuses (where available). Our sample includes single-resident 54 (46%), two-resident (18%), three or more-resident 55 56 homes (4%), and a not-reported category (32%). Residents can be categorized as young (age <35, 14%), 57

111

112

113

middle-aged (age 35-64, 9%), senior (age >64, 65%), and 58 a not-reported category (12%). Education levels of 59 60 residents in our dataset varies, including a high school 61 diploma (10%), bachelor's degree (19%), master's level (20%), doctorate degree (15%), and a not-reported 62 category (36%). Our entire 99 smart home dataset 63 includes people who are healthy (57%) and those with 64 targeted health ailments (23%) such as mild cognitive 65 66 impairment (MCI) (9%), as well as those whose health status was not reported (20%). 67

While we have collected a large set of sensor data for 68 this analysis, the data may not be representative of the 69 population as a whole. Thus, we employ different indices 70 to determine how representative our data are of the 71 72 national population. Information statistics (Shannon index) and dominance (Simpson index) indices are 73 74 utilized to identify and quantify both the richness (number of subgroups present) and abundance (the 75 number of individuals per subgroup) of our smart home 76 dataset in comparison with the US population. We also 77 utilize mean and variance analysis to investigate the 78 composition of the dataset. Further, we utilize Jaccard's 79 coefficient index, a value between 0 (not similar) and 1 80 (identical), to compare the similarities between our 81 smart home dataset and the US population. The data of 82 83 the US population in 2010 is collected from the census 84 government website [60]-[63].

85 For the information statistics, Fig. 4 shows that the value of Shannon indexes for age in our sample (1.01) is 86 close to that at the national average (1.03), reflecting the 87 richness and abundance of our smart home dataset. We 88 further use the Simpson index to analyze the dominant 89 subgroup (see Fig. 5) as well as mean and standard 90 deviation to analyze the composition in both datasets 91 (see Fig. 6). In Fig. 5, for the category of education level, 92 93 both our smart home dataset and the national population 94 have a Simpson index greater than 0.6, reflecting diverse education levels and no dominant subgroup. For age and 95 number of residents, the Simpson index in our smart 96 home dataset is less than 0.5 (dominant subgroups may 97 exist). The mean and deviation charts (see Fig. 6) show 98 that the mean value of age in our sample (76) is over 99 twice that at the national level (37). In respect to 100 101 household size, our sample does include fewer residents on average (1.4) than that in the national population 102 103 (2.6). Our dataset is more representative of senior 104 residents and low-population homes. The values of the Jaccard index for the categories of age(s), number of 105 residents, and education level are 0.14, 0.14, and 1, 106 respectively. That is, in the category of education level, 107 our dataset is similar as the national population. A 108 complete list of home descriptive parameters is provided 109 in the supplementary material. 110



Figure 4. Shannon Index (information statistics index) of 11 smart home dataset and national population in 2010. 12



24 Figure 5. Simpson Index (dominance index) of smart home dataset and national population in 2010. 25 26

27 To analyze differences in behavior models between 28 different population subgroups for individual activities, 29 we select two subgroups among our smart home participants who are matched in terms of age and 30 number of residents. The first subgroup consists of 31 senior residents who are healthy and living alone 32 (Subgroup H). This represents a baseline group for a 33 comparison to senior residents who are living alone and 34 have chronic health ailments. Most of these residents had 35 multiple health problems. The most significant limiting 36 conditions included mild cognitive impairment (N = 4)37 and mild dementia (N = 3). Further, one resident has 38 39 Parkinson's Disease, 4 people have mobility limitations, 2 residents have lung problems (chronic 40 hypoxia/chronic obstructive pulmonary disorder), 2 41 participants have atrial fibrillation, and 1 resident has 42 macular degeneration. . To keep our sample sizes as large 43 as possible, we did not constrain education level for 44 these participants. There are 16 homes included in 45 Subgroup H and 17 homes included in Subgroup NH. To 46 study an entire sequence of activities, we selected a 47 48 home with over five years of collected data. The resident 49 of this selected home also experienced health changes during the data collection period. Specifically, the 50 resident has vision and mobility problems. 51

- 52
- 53
- 54
- 55

Dataset Representativeness Analysis 56 the number of residents 90 80 M 70 д бо N 50 60 60 V 40 L 30 24 21 E 20 Î details 63 10 1.4 2,6 the number of residents age(s) education level national population in 2010 the entire smart homes dataset

67 Figure 6. Mean and variance of each category from smart home dataset and national population in 2010. 68

3.2. Data Cleaning 69

57

58

59

61

62

64

65

66

70 Before we fit a model to our smart home data, we preprocess the inter-arrival times of activity segments to 71 72 remove noise. Noise can arise in smart home data due to 73 issues including sensor failure, visitors in the home, activity recognition/segmentation errors, or changes 74 that are made to the environment. While some of the 75 outliers may represent behavioral changes that need to 76 be captured, others represent errors in the data 77 collection process and are best removed. 78

79 Our study focuses on large sets of continuous realvalued data. Longitudinal data collected from real homes 80 is also subject to noise resulting from imperfect sensors 81 82 and related system issues. As a result, we first apply outlier detection to the activity inter-arrival times by 83 using a threshold exceedances approach, a principal 84 approach found in extreme value theory (EVT). This 85 approach allows us to test a range of threshold values u 86 and identify outliers which have values above the 87 threshold. For each candidate threshold, we fit a 88 distribution for the excesses, the difference between the 89 outlier and the threshold. Lower thresholds tend to bias 90 the excess model by categorizing a large amount of data 91 as outliers. Higher thresholds lead to a greater variance 92 of the excess distribution because of the small number of 93 outliers. The standard rule is to choose a threshold as 94 low as possible so that the excesses fit a reasonable 95 96 distribution [64], [65].

Here, a reasonable distribution is governed by two 97 98 factors. First, we strive for a balance between bias and variance of the excesses distribution. Mean residual life 99 plots help visualize this balance. A mean residual life plot 100 graphs the mean value of excesses as a function of the 101 threshold value. We select the threshold value at the 102 lowest threshold value to show linearity in the plot. 103 Linearity indicates that the bias and variance of the 104 excess distribution are nearly evenly balanced. As an 105 example of this approach, Fig. 7 hows the mean residual 106 life plot with 95% confidence intervals for the inter-107 arrival times of the "Personal Hygiene" activity for our 108

Using Continuous Sensor Data to Formalize a Model of In-Home Activity Patterns

sample of 99 homes. In Fig. 7(a), the x-axis shows a range
 of threshold values (*u*) from the minimum to the
 maximum values observed Personal Hygiene inter arrival times. The y-axis shows the mean excess (the
 mean of the excess times above the threshold) for each
 threshold value.



Figure 7. (a) The mean residual life plot of Personal 30 Hygiene inter-arrival times for the entire dataset. Dashed 31 lines represent the 95% confidence interval. The solid 32 33 line plots the threshold value against the mean excess. Both the threshold and excess values are in hours. The 34 graph is near-linear at u = 1000 hours; (b) The mean 35 residual life plot of the Personal Hygiene inter-arrival 36 time for a range of thresholds from 5 to 20. Dashed lines 37 38 represent the 95% confidence interval. The solid line represents the threshold value against the mean of 39 40 excesses. This graph is near-linear at u = 17 hours. 41

To balance the bias and variance of the excess 42 distribution, we identify a threshold value where the 43 44 solid line in Fig. 7(a) shows linearity. We notice that the graph appears to be approximately linear around u =45 1000 hours. Since the value at the 99th percentile of the 46 entire dataset of the inter-arrival time of Personal 47 Hygiene is 13.1 hours, the value 1000 (hours) can be 48 49 considered a high threshold. To reduce the variance of excess model fitting, we choose a threshold based on a 50 mean residual life plot for a range of thresholds from 5 to 51 52 20 (hours) as shown in Fig. 7(b). This plot suggests an upper threshold (the maximum inter-arrival times) of 53 approximately 17 hours with 2835 outliers out of 54 366,441 data points, or 0.774% of the total number of 55 data points. 56

57 In the second step, we refine the threshold choice to 58 ensure that the shape parameter (affecting the shape of

a distribution instead of shifting or stretching/shrinking 59 the distribution) and modified scale parameter (affecting 60 the stretching/shrinking of a distribution) of the excess 61 distribution are quantifiably stable. To do this, we select 62 the lowest threshold value, near the approximated 63 64 threshold from the first step, for which both the 65 estimated shape parameter remains near-constant and the estimated modified scale parameter is near-linear 66 67 [64].



85 Figure 8. Parameter estimates against a range of thresholds from 10 to 35 hours from the Personal 86 Hygiene inter-arrival times. We select the lowest value (u 87 18 hours, the maximum inter-arrival times) of 88 = thresholds, near the approximated threshold (u = 1789 90 hours, the maximum inter-arrival times), for which the estimated shape (affecting the shape of the distribution 91 rather than shifting or stretching/shrinking it) 92 93 parameter remains near-constant and the estimated 94 modified scale (affecting the stretching/shrinking of a distribution) parameter is near-linear. 95 96

97 Using the same example of Personal Hygiene inter-98 arrival times as in the first approach, Fig. 8shows shape and re-parametrized scale parameters as a function of 99 100 alternative threshold values. Based on the plots in Fig. 8, which offer a model-based analysis of excess, we choose 101 a threshold of u = 18 hours (the maximum inter-arrival 102 times) with 2730 outliers out of 366,441 data points 103 (0.745%). The perturbations among the excesses are 104 105 small relative to sampling errors based on the stability of 106 parameter estimates in Fig. 8.

107 4. MODEL FITTING

Modeling human behavior from smart home sensors
provides a unique perspective not only to investigate
behavioral norms but also to compare these norms
between population subgroups. Once differences in
patterns are discovered, they can be used to better
understand the impact of personal characteristics, such





as age, health conditions, and education on daily
routines. They can also be used to automate diagnoses
and predict additional behavioral features of individuals
within a group. In this section, we provide details of our
model fitting procedures. To illustrate the process, we
focus on the Personal Hygiene activity observed from all
the sampled smart homes.

20 We use the data below the selected EVT threshold to determine which distribution best describes the sensor-21 based activity data. We model data histograms using 82 22 different probability distribution functions. In our study, 23 24 we employ the Freedman-Diaconis (F-D) rule to select 25 the size of the bins for the histogram. In general, three well-known rules are used to calculate bin widths: the 26 Sturges, Scott, and F-D rules. The Sturges rule is 27 applicable when the data is from symmetric and 28 Gaussian distributions [66]. The Scott rule works well for 29 non-Gaussian distributions but refers to sample sizes 30 between 50 and 500. For larger samples as in our study, 31 the Scott rule estimates a smaller number of histogram 32 bins, leading to over-smoothing. Over-smoothed 33 34 histograms provide limited information on the shape of 35 the underlying distribution [67]. The F-D rule is a robust method that substitutes the estimated standard 36 37 deviation in the Scott rule with a multiple of the interquartile range. Thus, the F-D rule ensures 35% 38 more bins than from the Scott rule as well as keeps the 39 property of the Scott rule for non-Gaussian distributions 40 [68], [69]. 41

42 Next, we use the smallest summation of the square errors (SSE) value to compare distributions and select 43 44 the distributions that best fit the data. SSE is calculated 45 by determining the difference between the data and the 46 fitted distribution. Here, we give an example of the distribution fitting and selection process for the Personal 47 Hygiene inter-arrival times from all smart homes. The 48 results of modeling the remaining activities and 49 comparing the two subgroups are provided as 50 supplementary material. While we use 82 distributions 51 to fit the histogram, Fig. 9 shows the top 15 fitted 52 distributions for the Personal Hygiene inter-arrival 53 54 times from all the smart homes. The Pareto distribution 55 is selected as one of these top distributions. Fig. 10

Figure 10. The SSE values for the top 15 fitteddistributions. Smaller SSE values indicate a better fit.

summarizes the SSEs between the smart home data and 59 the top 15 fitted distributions. The Pareto distribution 60 has the same order-of-magnitude errors (10-1) as the 61 62 other top distributions. We hypothesize that the Pareto distribution provides a close approximation to the top 63 fitted distribution for Personal Hygiene inter-arrival 64 times based on the sampled 99 smart homes. We test this 65 hypothesis by both visualizing the fitting and 66 determining the significance of the difference in fit 67 between the Pareto distribution and the top-fitting 68 distribution. 69

70 First, the figure on the left side of Fig. 11 shows the 71 fit of the Pareto distribution and the top-fitting 72 distribution for the Personal Hygiene inter-arrival times across all smart homes. Figures on the right side, (b) and 73 (c), of Fig. 11, respectively, show portions of the fitted 74 Pareto distribution for the Personal Hygiene inter-75 arrival times from 0 to 6.5 (hours) and from 3 to 18 76 77 (hours). Based on Fig. 11, the fitted Pareto distribution well approximates the shape of the histogram of the 78 Personal Hygiene inter-arrival time in the entire dataset, 79



92 Figure 11. (a) Distribution fitting between the Pareto distribution and the top-fitting distribution. Bins are 93 indicated by the x-axis. The y-axis represents frequency, 94 the amount of data included in each bin divided by the 95 total amount of data. The SSE of this Pareto distribution 96 97 is 0.3. (b) Portion of the graph corresponding to shorter inter-arrival times (from 0 to 6.5 hours). (c) Portion of 98 the graph corresponding to longer inter-arrival times 99 (from 3 to 18 hours). 100

though the distribution did not capture everything from
 the histogram. For example, a hump exists (see Fig. 11

3 (c)) around hours 8 through 10 with frequencies 0.01 to

4 0.02. Because, we are capturing a general view of indoor

- 5 behavior patterns rather than modeling each detail of a
- 6 single activity, we may also be observing overfit.

Second, to validate that the Pareto model provides a 7 statistically significantly-similar fit to the top-fitting 8 9 distribution, we perform a t-test analysis with the null hypothesis that the two distributions have identical fit 10 scores. Given the activity inter-arrival times and the 11 estimated distribution parameters, we generate the 12 values of probability density functions for the Pareto 13 distribution and the top-fitted distribution. Next, we split 14 the two sets of values into 60 subsets and perform a 15 16 paired t-test on the means for each subset.

17 For Personal Hygiene inter-arrival times across all 18 sampled homes, the p value is 0.153 with the t-statistic -19 1.443. For the null hypothesis that the two distributions have identical average scores, a small p value (<= 0.05) 20 leads to rejecting the null hypothesis and a large p value 21 (>= 0.05) leads to accepting the null hypothesis. Thus, we 22 accept the null hypothesis, and the Pareto distribution 23 can be considered approximately as strong as the top-24 fitting distribution for the Personal Hygiene inter-arrival 25 26 times from the entire collection of smart homes. Similar 27 results were observed for all selected seven activities. 28 Based on both the visualization and t-test, the Pareto 29 distribution provides a strong fit for this activity data and the properties of the Pareto distribution, for example the 30

- 31 80/20 rule, provide opportunities for future analyses
- 32 and investigations of hypotheses that model indoor33 behavior patterns.

34 5. RESULTS

To understand the general principles behind humanbehavior in everyday environments and to compare the

behavioral norms between population groups, we 37 perform the same procedures described in Sections 3 38 and 4 for each recognized activity both across all 99 39 smart homes and among two older adult subgroups 40 (Subgroup_H and Subgroup_NH). In addition, using the 41 42 same procedures as in Sections 3 and 4, we study a 43 holistic behavior routine in one home as a combination of all detected activities. 44

45 For the data from 99 smart homes, before performing outlier detection on the inter-arrival times of 46 47 each activity, we visualize some statistics to gain an intuitive understanding of the data (see Fig. 12). In these 48 graphs, we observe that the maximum value of the inter-49 arrival times of each activity is relatively large (>= 103 50 hours). There are multiple possible explanations for 51 52 these large values, including sensor failures and the 53 resident's absence from the home during travel. The 99th 54 percentile of inter-arrival times for each activity is in the 55 range of 101 to 102 hours. That is, approximately 99% of occurrences for each activity exhibit small inter-arrival 56 times, thus only the top 1% of inter-arrival times 57 demonstrate these large values of interest. The mean 58 inter-arrival time for each activity is in the range of 100 59 to 101 hours. For example, the mean value of the inter-60 arrival times of two successive Cook activities from 99 61 62 smart homes is 3.5 hours, which gives us a generalized 63 view about the gap between two successive Cook 64 activities. In Fig. 13, we notice that over 99.5% of the data 65 are kept after removing outliers. In addition, the threshold value of each activity is above their mean value 66 (except activity Eat) and sometimes above the 99th %. 67

After filtering the outliers (using methods described
in Section 3), we perform model fitting as described in
Section 4. The summarized result of fitting of activity
inter-arrival times is shown in Fig. 14 and Table 1. In Fig.
14, we noticed that the shape parameter of the Pareto
distribution for Sleep inter-arrival times from the entire
smart home dataset is less than one. This means that the



Figure 12. A summary of inter-arrival times of all activities for all smart homes before performing outlier detection. The
 results include the mean value of the dataset (mean), the 99th percentile (99th%) and maximum value of the dataset
 (max).

Using Continuous Sensor Data to Formalize a Model of In-Home Activity Patterns



Figure 13. Summarized results of all activities for all smart homes. The results include the selected upper threshold (*u*), the percentage of data that remains after removing the outliers (*r*), and the number of outliers (*o*).



Figure 14. The Pareto distribution of each activity givensimulated inter-arrival times from 0 to 40 hours.

expected start time of a Sleep activity relative to the 32 previous Sleep occurrence approaches infinity. This 33 result occurs because the mean value of the dataset is 34 influenced by the largest single value. This may occur in 35 finite-size samples when an outlier causes the mean to 36 become arbitrarily large. The Sleep activity arrival times 37 cannot be adequately captured 38 therefore bv 39 distributions and thus we will use quantiles to describe the data spread of the Sleep activity. 40

Based on Fig. 15, we observe that for each activity, the SSE of the Pareto distribution fit is relatively small, in the range of 10-3 to 10-1. Furthermore, the Pareto distribution for each activity is approximately as strong as the top-fitting distribution based on the large t-test p values (>= 0.05) in Table 1 (the null hypothesis is that the 56 two distributions have identical mean scores).

To further interpret the behavioral norms, we compare the selected thresholds (*u*) and the Pareto shape parameters (α) for the entire sampled population and among population subgroups (see Figures 16 and 17). In Fig. 16 (a), the threshold of the inter-arrival time of Work, Relax, Cook, Eat, and Sleep in Subgroup NH is smaller than the corresponding threshold in Subgroup H.

One possible explanation for this difference is that 64 individuals in Subgroup NH may be unable to sustain 65 66 long periods of one activity, thus creating bursts of 67 activities with short breaks. The phenomena might be 68 due to physical health ailments, such as mobility or 69 stamina difficulties that may require periods of rest. In 70 addition, participants with cognitive limitations may be 71 more likely to get distracted or experience difficulty 72 remembering to quickly return to a task following an activity interruption, resulting in the need to reinitiate an 73 activity within a short period of time. 74

Fig. 17 ows the values of the shape parameters of 75 76 the Pareto distributions for the individual selected 77 activities, the variations of which may be due to variations in health conditions. In Fig. 17(a), we observe 78 79 that the shape parameters for activities Personal Hygiene, Relax, Cook, and Sleep in Subgroup H are 80 smaller than the shape parameters in Subgroup NH. The 81 smaller the shape parameter, the heavier the tail is of the 82 Pareto distribution. That is, longer starting times 83 between two successive activities occur more frequently 84 in Subgroup H, possibly due to fewer interruptions. 85 86

47 48

31

1

2

3

4 5

6

7 8

9

10

11

12

13

14 15

10	Table 1. The p value of the t-test between the top- fitting distribution and the rate of distribution (p) .								
49		Тор	Fitted	р		Тор	Fitted	p value	
50		Distribution	Vä	alue	D	istribution			
51	Work	Burr		0.540	Relax	Inverse Ga	ussian	0.283	
52	Wash Dishes	Fatigue life		0.312	Cook	Lévy		0.247	
53	Personal Hygiene	Fatigue life		0.153	Eat	Fatigue life	e	0.095	
54									

Table 1. The purplue of the t test between the ten, fitting distribution and the Darote distribution (n)

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

76

77

80

106



Figure 15. The SSE values of the top-fitting distribution 10 and the fitted Pareto distribution. 11

12

Based on the above observations, we hypothesize 13 that activity inter-arrival times may be used to predict 14 subgroup classifications. We also notice that the Work 15 16 activity (typically working at a desk or on a computer in 17 an office area of the home) exhibits a large difference 18 (0.36) in model shape parameters between the 19 subgroups. To predict subgroup classifications, we currently use Work inter-arrival times from both 20 subgroups and then utilize a random forest algorithm 21 with 10-fold cross validation. The average accuracy of 22 predictions is 0.814 and the standard deviation is 9.5. 23 The precision of predicting subgroup classification is 24 0.784 and the recall is 0.740. 25

26 Further, to validate that the random forest algorithm 27 provides a statistically significantly-better prediction than that from random guesses, we perform a t-test 28 29 analysis with the null hypothesis that the mean difference of the prediction results from the two 30 algorithms are equal to zero. The p value of the t-test is 31 0.0001 with the t-statistic 10.43. Since a small p value (<= 32 0.05) leads to rejecting the null hypothesis, we concludes 33 that the prediction results from the random forest 34 classifier using model parameter attributes are 35 statistically significantly-better than those from random 36 guesses. These results indicate that the formal model 37 does indeed reflect differences in behavior patterns or 38 39 population subgroups and can help us understand behavioral impacts of traits such as chronic health 40 conditions. 41

In addition, the shape parameters for activities 42 Wash Dishes, Personal Hygiene, Cook, and Eat in the 43 combined dataset are larger than parameters for either 44 of the subgroups. That is, the shorter starting times 45 between two successive activities occur more frequently 46 in the combined dataset. This may be due to the number 47 48 and diversity of residents in the combined dataset. 49 Homes with multiple residents, young, or middle-age residents have a higher frequency of shorter inter-arrival 50 times than that for single senior residents. We also 51 noticed that the shape parameter of Work in the 52 53 combined dataset is smaller than either of the subgroups 54 (larger gaps between two successive Work activity occurrences in the combined dataset exist), possibly 55



Figure 16. Summarized results of upper thresholds for 75 ADLs in the complete dataset and for Subgroups H and NH. (a) Upper thresholds as a function of activity category. (b) Upper thresholds as a function of the 78 subgroup. 79



Figure 17. Summarized results of the Pareto shape 101 parameters for routine activities in the complete dataset 102 103 and for Subgroups H and NH. (a) Pareto shape 104 parameters as a function of activity category. (b) Pareto 105 shape parameters as a function of the subgroup.

107 due to the fluctuation of residents' schedules, such as when the residents' are travelling, while seniors often 108 109 have more stable schedules.

In Fig. 17(b), we observe that in the combined 1 56 2 dataset, the activities Relax and Eat have approximately 57 the same value of the shape parameters (1.18 and 1.17, 3 58 respectively). In Subgroup NH, activities Wash Dishes 4 59 and Eat share the same value of the shape parameter 5 60 6 (1.08 for both). In Subgroup H, four activities, Personal 61 Hygiene, Relax, Cook, and Eat, have almost the same 7 62 values of the shape parameters (1.11, 1.11, 1.09, and 8 63 9 1.11, respectively). Given the similar shape parameters, 64 we propose to utilize bivariate or multivariate Pareto 10 65 distributions (its cumulative density functions are 11 shown in equations 1-3) to describe the combination of 12 67 activities in each group. That is, the interdependencies of 13 certain activities exist both at 99 smart homes and 14 among subgroups. For example, in Subgroup NH, we can 15 16 use a bivariate Pareto distribution to describe the 71 relationship between activities Wash Dishes and Eat. In 17 72 18 our study, all Pareto distributions are Pareto Type II. 73 19 The Summary of Quantiles of Sleep Inter-arrival Times 20

(a)

Hours

3rd Quantile

2nd Quantile/Median

21 22 23

24

25

26 27

28 29

30

31

35

SUBGROUP

1st Quantile



consequently in model parameters) for healthy and non-

healthy subpopulations given the activities we

Figure 18. Summarized results of Sleep inter-arrival times. (a) includes the first quantile (1st Qu.), the median
 value (Median) / 2nd quantile (Median), the third quantile (3rd Qu.). (b) includes the minimum value of the
 dataset (Min.), the mean value (Mean), maximum value of the dataset (Max.).

36 The previous activities were tightly modeled as Pareto distributions. For Sleep inter-arrival times, the 37 shape parameters are less than one for the entire sample 38 39 of 99 smart homes and among population subgroups (see Fig. 17(a)). Statistically, this implies that the 40 expected inter-arrival time approaches infinity. This 41 result occurs because the mean value of the dataset is 42 influenced by the largest value of the dataset. For a 43 dataset of finite size, the sample has a finite value and so 44 does the mean. But the more samples we have, the larger 45 value of the mean. That is, the estimate of the mean is 46 divergent when the size of the dataset goes to infinity. 47 Since we cannot find a fitted distribution to adequately 48 49 describe the pattern of the Sleep data below the selected thresholds, we utilize quantiles to understand the data 50 spread (see Fig. 18). We notice that all the values in 51 Subgroup H are greater than those in Subgroup NH. This 52 is likely because residents with health ailments tend to 53 54 experience more interrupted sleep. As this discussion highlighted, we do see differences in behaviour (and 55

90 Studying activity classes separately cannot provide a comprehensive view of a person's entire routine. To 91 understand all activities comprising a routine, we select 92 one home to investigate patterns of all activities. We 93 combine both the predefined (and labelled) seven 94 95 activities and the remaining clustered activities. The appropriate number of clusters is selected when no 96 significant change of the sum of squared distances occurs 97 in the elbow curve. That is, the optimal number of 98 clusters is near the elbow part of the curve. Based on Fig. 99 100 19, we select *k* = 10. The resident is a single senior whose health status transitioned from healthy to having vision 101 and mobility problems during the course of data 102 collection. The time period of the experiment for this 103 home is 65 months (from 2011-06-14 to 2016-11-10). 104

We perform the same process of data processing and model fitting as described in Sections 3 and 4. Fig. 20 shows that the Pareto distribution also fits the interarrival times for all routine activities. The threshold and shape parameter of the inter-arrival times of all activities are 5.8 hours and 1.17, respectively. The sum of square
 error for the fitted Pareto distribution is 1.8. Given the
 values of the shape (1.17) and scale (0.29) parameters,
 we confirm that 17% of the total number of inter-arrival
 points, which occurs in the tail, comprises 80% of the
 total inter-arrival hours, and 20% of the total number of

7 inter-arrival points, comprises 74% of the total inter-8 arrival hours.

9 To further investigate the relationship between the fat tail of the Pareto distribution and the resident's health 10 status, we first look at the 20% of inter-arrival times that 11 occur in the tail and represent the large gap between two 12 successive activities, and then manually examine the 13 sensor data that corresponds to these gaps. Our 14 investigation is summarized in Fig. 21. We notice that the 15 16 successive activities, Sleep and Bed-Toilet Transition, 17 occur in 40% of the cases lying in the tail. We hypothesize 18

that the large gap and high frequency of these two 55 activities are symptomatic of the resident's health 56 problems. We validate the hypothesis by comparing the 57 provided health information with the sensor-based night 58 time walking duration (minutes) for the corresponding 59 60 dates (see Fig. 22). Average night time walking duration is calculated based on the time that elapses between the 61 end of a sleep activity and the beginning of the following 62 63 bed toilet transition, given that the distance between bed and bathroom is constant and night-time bathroom trips 64 typically involve direct routes. 65

We observe an increase in the average walking duration from August 2014 to November 2014. Selfreported mobile difficulty also increases during that time, provide a possible explanation for this change. We also notice that from December 2014 to March 2015 the sensor data reflects an decrease in the average walking duration, while the self-reported mobility difficulty



Figure 21. Spread of the two successive activities in 20% of the total inter-arrival times that occur in the tail of the distribution.

109

53

54



Figure 22. Compare the monthly average walking speed captured by sensor data with the self-reported mobility difficulty.

consistently drops from 3 to 1 (on a scale from 1= no
difficulty to 5= tremendous difficult). The results
provide evidence that the large time gaps and high
frequency of Sleep and Bed-Toilet Transition activities in
this particular home are related to the resident's health
status.

24 6. DISCUSSION, LIMITATIONS, AND 25 DIRECTIONS FOR FUTURE 26 RESEARCH

27 In this paper, we propose formal methods for 28 modeling human routines in everyday environments. We found that the Pareto distribution fits many activities, 29 thus providing unique insights into behavior norms for 30 the entire sampled population and behavior variations 31 between population subgroups. Further, we discover 32 that several activities in certain groups can be described 33 by multivariate Pareto distributions. We also explore the 34 35 pattern of all activities as a routine in one home and its 36 relationship with the resident health ailments.

37 When applying our analysis to smart home data, we 38 find that activities follow a non-Poisson process and the inter-arrival times of individual activities as well as all 39 activities within a holistic routine fit a Pareto 40 distribution. The findings may provide useful 41 information to further investigate potential behavior 42 changes that might be related to health problems. 43 Limitations of this study include sensitivity of the models 44 45 to noise in the sensor data, addressed in part by the 46 outlier filtering process. In addition, the Pareto 47 distribution fits the data well but does not fully describe 48 all routine details. For example, the small hump in the histogram of the data (see Fig. 11 (c)) exists with 49 frequency in a range of 0.01 to 0.02. A mixture model 50 may be introduced to capture the hump for greater 51

70 model detail. In Section 3, the selection of the threshold may impact the parameters of the Pareto distribution, 71 especially when investigating the distribution's 20% tail. 72 However, since we look at the highest percentage of two 73 successive activities that lie in the tail, the impact of the 74 similar threshold/shape parameters may be small. 75 Further, the current two subgroups only consider single 76 77 residents instead of multiple residents, though the 78 experiments that evaluate the entire 99 smart homes do include multiple residents with diverse backgrounds. 79 80 The problem of tracking, recognizing, and analyzing multi-resident behavior is an ongoing challenge, 81 although Wang et al. discuss one possible strategy for 82 83 multi-resident tracking in smart homes [70].

In addition, one can observe that our initial model 84 oversimplifies the complexity of human behavior. For 85 the purpose of this present study, we intentionally kept 86 87 the model simple and focused on automatically-detected activity timings in smart environments. However, 88 89 development of more sophisticated models combining 90 other parameters including social interactions, circadian rhythms, night time relative walking speed, and 91 movement patterns in and out of the home can further 92 boost our ability to understand and reproduce the 93 structure of human activities. 94

Additionally, further analysis of each activity and all 95 activities in a routine will allow us to address specific 96 97 questions that have been asked in the literature. For example, we could provide evidence to support or deny 98 99 the hypothesis that human behavior in certain groups is 100 random or Markovian [71]-[74]. We could also examine whether the 80/20 rule in which the largest-distance 101 movements (20% of movement distances) occur with 102 80% of total inter-arrival hours applies in home 103 104 environments. Finally, future work can quantify the 105 predictability of behavior parameters for different population groups and types of sensed data. 106

1 ACKNOWLEDGMENTS

- 2 This material is based upon work supported by the
- 3 National Science Foundation under Grant No. 1543656.
- 4 The authors would like to thank Samaneh
- 5 Aminikhanghahi for her help with change point analysis.6

7 **REFERENCES**

- 8 [1] G. Last and M. Penrose, Lectures on the Poisson 9 process, vol. 7. Cambridge University Press, 2017. [2] F. Jovan, J. Wyatt, N. Hawes, and T. Krajník, "A 10 Poisson-spectral model for modelling temporal 11 patterns in human data observed by a robot," in 12 Intelligent Robots and Systems (IROS), 2016 13 14 IEEE/RSJ International Conference on, 2016, pp. 4013-4018. 15
- 16 [3] R. G. Gallager, "Poisson processes," *Stoch.*17 *Process. Theory Appl.*, pp. 74–108, 2013.
- [4] A. Ihler, J. Hutchins, and P. Smyth, "Adaptive event detection with time-varying poisson processes," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2006, pp. 207–216.
- 23 [5] J. G. Oliveira and A.-L. Barabási, "Human
 24 dynamics: Darwin and Einstein correspondence
 25 patterns," *Nature*, vol. 437, no. 7063, p. 1251,
 26 2005.
- A. Vázquez, J. G. Oliveira, Z. Dezsö, K. Il Goh, I.
 Kondor, and A. L. Barabási, "Modeling bursts and heavy tails in human dynamics," *Phys. Rev. E* -*Stat. Nonlinear, Soft Matter Phys.*, 2006.
- [7] R. F. Grais, J. H. Ellis, and G. E. Glass, "Assessing
 the impact of airline travel on the geographic
 spread of pandemic influenza," *Eur. J. Epidemiol.*,
 vol. 18, no. 11, pp. 1065–1072, 2003.
- A. Pieropan, C. H. Ek, and H. Kjellström,
 "Functional object descriptors for human
 activity modeling," in *Robotics and Automation*(ICRA), 2013 IEEE International Conference on,
 2013, pp. 1282–1289.
- 40 [9] O. Kwon, W.-S. Son, and W.-S. Jung, "The double
 41 power law in human collaboration behavior:
 42 The case of Wikipedia," *Phys. A Stat. Mech. its*43 *Appl.*, vol. 461, pp. 85–91, 2016.
- L. L. Constantine, "Human activity modeling: toward a pragmatic integration of activity theory and usage-centered design," in *Humancentered software engineering*, Springer, 2009, pp. 27–51.
- 49 [11] G. Gay and H. Hembrooke, Activity-centered
 50 design: An ecological approach to designing
 51 smart tools and usable systems. Mit Press, 2004.
- L. L. Constantine and L. A. D. Lockwood, Software
 for use: a practical guide to the models and
 methods of usage-centered design. Pearson
 Education, 1999.
- A. Bees, N. York, and A. Barabasi, "The origin of
 bursts and heavy tails in human dynamics," *Nature*, vol. 435, no. 7039, pp. 207-211, 2005.

- J. Leskovec, M. McGlohon, C. Faloutsos, N. Glance, and M. Hurst, "Patterns of cascading behavior in large blog graphs," in *Proceedings of the 2007 SIAM international conference on data mining*, 2007, pp. 551–556.
- [15] J. Ratkiewicz, S. Fortunato, A. Flammini, F.
 Menczer, and A. Vespignani, "Characterizing and modeling the dynamics of online popularity," *Phys. Rev. Lett.*, vol. 105, no. 15, p. 158701, 2010.
- [16] R. Kumar, M. Mahdian, and M. McGlohon,
 "Dynamics of conversations," in *Proceedings of*the 16th ACM SIGKDD international conference
 on Knowledge discovery and data mining, 2010,
 pp. 553–562.
- 73 [17] H. Li *et al.*, "Characterizing smartphone usage patterns from millions of android users," in
 75 Proceedings of the 2015 Internet Measurement Conference, 2015, pp. 459–472.
- Y. Gandica, J. Carvalho, F. S. Dos Aidos, R.
 Lambiotte, and T. Carletti, "Stationarity of the inter-event power-law distributions," *PLoS One*, vol. 12, no. 3, p. e0174509, 2017.
- 81 [19] I. Tsompanidis, A. H. Zahran, and C. J. Sreenan,
 "Mobile network traffic: A user behaviour
 model," in 2014 7th IFIP Wireless and Mobile
 Networking Conference (WMNC), 2014, pp. 1–8.
- 85 [20] L. Yu, P. Cui, C. Song, T. Zhang, and S. Yang, "A 86 temporally heterogeneous survival framework 87 with application to social behavior dynamics," in 88 Proceedings of the 23rd ACM SIGKDD 89 International Conference on Knowledge 90 *Discovery and Data Mining*, 2017, pp. 1295–1304.
- 91 [21] T. M. Scholz, "The human role within
 92 organizational change: A complex system
 93 perspective," in *Change management and the*94 *human factor*, Springer, 2015, pp. 19–31.
- 95 [22] P. Andriani and B. McKelvey, "Perspective—
 96 From Gaussian to Paretian thinking: Causes and
 97 implications of power laws in organizations,"
 98 Organ. Sci., vol. 20, no. 6, pp. 1053–1071, 2009.
- 99 [23] Y. U. Saito, T. Watanabe, and M. Iwamura, "Do larger firms have more interfirm relationships?,"
 101 *Phys. A Stat. Mech. its Appl.*, vol. 383, no. 1, pp. 102 158–163, 2007.
- 103 [24] M. C. Gonzalez, C. A. Hidalgo, and A.-L. Barabasi,
 "Understanding individual human mobility
 patterns," *Nature*, vol. 453, no. 7196, pp. 779–
 106 782, 2008.
- 107 [25] I. Rhee, M. Shin, S. Hong, K. Lee, S. J. Kim, and S.
 108 Chong, "On the levy-walk nature of human
 109 mobility," *IEEE/ACM Trans. Netw.*, vol. 19, no. 3,
 110 pp. 630–643, 2011.
- 111
 [26]
 C. Song, Z. Qu, N. Blumm, and A.-L. Barabási,

 112
 "Limits of predictability in human mobility,"

 113
 Science (80-.)., vol. 327, no. 5968, pp. 1018–

 114
 1021, 2010.
- 115 [27] W.-Y. Zhu, W.-C. Peng, L.-J. Chen, K. Zheng, and X.
 116 Zhou, "Modeling user mobility for location promotion in location-based social networks," in
 118 Proceedings of the 21th ACM SIGKDD

1		International Conference on Knowledge	61		r
2		Discovery and Data Mining 2015 nn 1573-1582	62		r
2	[20]	S Hong Human movement natterns mobility	62	[42]	P
ر ۲	[20]	s. Hong, Human movement patterns, mobility	03	[43]	Г ((
4		North Carolina State University 2010	64 CF		
5	[20]	North Carolina State University, 2010.	65		S
6	[29]	K. Zhao, M. Musolesi, P. Hui, W. Rao, and S.	66	F 4 43	V
7		Tarkoma, "Explaining the power-law	67	[44]	J
8		distribution of human mobility through	68		S
9		transportation modality decomposition," Sci.	69		r
10		<i>Rep.</i> , vol. 5, p. 9136, 2015.	70		l
11	[30]	R. Gallotti, A. Bazzani, S. Rambaldi, and M.	71		2
12		Barthelemy, "A stochastic model of randomly	72	[45]	(
13		accelerated walkers for human mobility." Nat.	73		Y
14		<i>Commun.</i> , vol. 7, p. 12600, 2016.	74		a
15	[31]	T Kurashima T Althoff and I Leskovec	75		2
16	[31]	"Modeling interdependent and periodic real	75		r
10		would action acquences" in Drocoodings of the	70	[46]	T
17		2010 Marld Wide Web Conference on World Wide	77	[40]	J
18		2018 World Wide Web Conjerence on World Wide	/8		I
19		Web, 2018, pp. 803–812.	79		S
20	[32]	R. Rawassizadeh, E. Momeni, C. Dobbins, J.	80		1
21		Gharibshah, and M. Pazzani, "Scalable daily	81	[47]	F
22		human behavioral pattern mining from	82		h
23		multivariate temporal data," IEEE Trans. Knowl.	83		r
24		Data Eng., vol. 28, no. 11, pp. 3098–3112, 2016.	84		J
25	[33]	N. Banovic, T. Buzali, F. Chevalier, J. Mankoff, and	85		2
26		A. K. Dev. "Modeling and understanding human	86	[48]	Ι
27		routine behavior." in <i>Proceedings of the 2016 CHI</i>	87	L - J	a
28		Conference on Human Factors in Computing	88		2
20		Sustans 2016 nn 248-260	80		F
29	[24]	M Pan at al "Dissocting the Learning Curve of	00	[40]	L (
30	[34]	M. Fall et ul., Dissecting the Leaf hing curve of	90	[49]	C C
31		Duran diam of the 2010 CIAM Internetical	91		S
32		Proceedings of the 2019 SIAM International	92		e
33	5 a - 7	<i>Conference on Data Mining</i> , 2019, pp. 783–791.	93		C
34	[35]	H. Ghayvat, J. Liu, S. C. Mukhopadhyay, and X. Gui,	94	[50]	5
35		"Wellness sensor networks: A proposal and	95		"
36		implementation for smart home for assisted	96		а
37		living," IEEE Sens. J., vol. 15, no. 12, pp. 7341–	97		Ι
38		7348, 2015.	98	[51]	N
39	[36]	G. Laput, Y. Zhang, and C. Harrison, "Synthetic	99		b
40		sensors: Towards general-purpose sensing," in	100		2
41		Proceedings of the 2017 CHI Conference on	101	[52]	I
42		Human Factors in Computing Systems, 2017, pp.	102		C
43		3986–3999.	103		а
44	[37]	B. Lin <i>et al.</i> , "Analyzing the relationship between	104		C
45	[07]	human behavior and indoor air quality" I Sens	105		2
45		Actuator Networks vol 6 no 3 n 13 2017	105	[53]	6
40	[20]	A D Diagoras K E Deannis C Storgiou H Wang	100	[55]	E E
47	[30]	and P. P. Cunta "Efficient LoT based concer PIC	107		L N
48		Data callection and analysis in a	108		r T
49		Data collection-processing and analysis in	109	[] 4]	1
50		smart buildings," Futur. Gener. Comput. Syst., vol.	110	[54]	F
51		82, pp. 349–357, 2018.	111		ŀ
52	[39]	D. J. Cook, "Learning setting-generalized activity	112		а
53		models for smart spaces," IEEE Intell. Syst., vol.	113		ŀ
54		2010, no. 99, p. 1, 2010.	114		2
55	[40]	N. C. Krishnan and D. J. Cook, "Activity	115	[55]	Γ
56		recognition on streaming sensor data," Pervasive	116		ŀ
57		<i>Mob. Comput.</i> , vol. 10, pp. 138–154, 2014.	117		(
58	[41]	D. Cook and N. Krishnan, "Activity Learning from	118		6
59		Sensor Data." Wiley, 2014.	119	[56]	(
60	[42]	E. Kim, S. Helal, and D. Cook. "Human activity	120	r .1	C
	r - - 1	-,, 00011, 4011101			c

recognition and pattern discovery," *IEEE pervasive Comput.*, vol. 9, no. 1, pp. 48–53, 2009.

- A. Benmansour, A. Bouchachia, and M. Feham,
- "Multioccupant activity recognition in pervasive smart home environments," *ACM Comput. Surv.*, vol. 48, no. 3, p. 34, 2016.
- [44] J. Wan, M. J. O'grady, and G. M. O'hare, "Dynamic sensor event segmentation for real-time activity recognition in a smart home context," *Pers. Ubiquitous Comput.*, vol. 19, no. 2, pp. 287–301, 2015.
- [45] G. Fairchild, K. S. Hickmann, S. M. Mniszewski, S.
 Y. Del Valle, and J. M. Hyman, "Optimizing human activity patterns using global sensitivity analysis," *Comput. Math. Organ. Theory*, vol. 20, no. 4, pp. 394–416, 2014.
- y[46]J. Wang, Y. Chen, S. Hao, X. Peng, and L. Hu, "Deepalearning for sensor-based activity recognition: Absurvey," Pattern Recognit. Lett., vol. 119, pp. 3–c)11, 2019.
- [47] A. Helal, D. J. Cook, and M. Schmalz, "Smart home-based health platform for behavioral monitoring and alteration of diabetes patients," *J. Diabetes Sci. Technol.*, vol. 3, no. 1, pp. 141–148, 2009.
- [48] D. J. Cook, M. Schmitter-Edgecombe, L. Jönsson, and A. V Morant, "Technology-enabled assessment of functional health," *IEEE Rev. Biomed. Eng.*, vol. 12, pp. 319–332, 2018.
- [49] C. Chen, D. J. Cook, and A. S. Crandall, "The user side of sustainability: Modeling behavior and energy usage in the home," *Pervasive Mob. Comput.*, vol. 9, no. 1, pp. 161–175, 2013.
- [50] S. Aminikhanghahi, T. Wang, and D. J. Cook, "Real-time change point detection with application to smart home time series data," *IEEE Trans. Knowl. Data Eng.*, 2018.
- [51] M. Yuan, "Human dynamics in space and time: A brief history and a view forward," *Trans. GIS*, vol. 22, no. 4, pp. 900–912, 2018.
- [52] J. J. Davis and E. G. Conlon, "Identifying compensatory driving behavior among older adults using the situational avoidance questionnaire," *J. Safety Res.*, vol. 63, pp. 47–55, 2017.
- [53] C. Li, W. K. Cheung, J. Liu, and J. K. Ng, "Automatic Extraction of Behavioral Patterns for Elderly Mobility and Daily Routine Analysis," *ACM Trans. Intell. Syst. Technol.*, vol. 9, no. 5, p. 54, 2018.
- [54] A. F. Costa, Y. Yamaguchi, A. J. M. Traina, and C. Faloutsos, "Modeling temporal activity to detect anomalous behavior in social media," *ACM Trans. Knowl. Discov. from Data*, vol. 11, no. 4, p. 49, 2017.
- [55] D. J. Cook, A. S. Crandall, B. L. Thomas, and N. C. Krishnan, "CASAS: A smart home in a box," *Computer (Long. Beach. Calif).*, vol. 46, no. 7, pp. 62–69, 2013.
- (56)C. Truong, L. Oudre, and N. Vayatis, "A review of
change point detection methods," *arXiv Prepr.*

1 2	[57]	<i>arXiv1801.00718</i> , 2018. D. Picard, "Testing and estimating change-points
3 4		in time series, <i>Aav. Appl. Probab.</i> , vol. 17, no. 4, pp. 841–867. 1985.
5	[58]	S. Aminikhanghahi and D. J. Cook, "Using change
6		point detection to automate daily activity
7		segmentation," in <i>Pervasive Computing and</i>
8		Communications Workshops (Percom Workshops) 2017 IEEE International Conference
9 10		on 2017 nn 262–267
11	[59]	R. P. Adams and D. I. C. MacKay, "Bayesian Online
12	[]	Changepoint Detection," 2007.
13	[60]	"Census Age Information." [Online]. Available:
14		https://www.census.gov/data/tables/2010/de
15		mo/age-and-sex/2010-age-sex-
16	[(1]	composition.html.
17 18	[61]	"Census Disability Characteristics." [Online].
19		https://factfinder.census.gov/faces/tableservic
20		es/jsf/pages/productview.xhtml?pid=ACS 16 1
21		YR_S1810&prodType=table%0D.
22	[62]	"Census Educational Information." [Online].
23		Available:
24		https://www.census.gov/data/tables/2010/de
25		mo/educational-attainment/cps-detailed-
26	[(0]	tables.html%0D.
27	[63]	"Census Households Information." [Unline].
28		Available:
29		series / demo / families / households html%0D
31	[64]	S. Coles, I. Bawa, L. Trenner, and P. Dorazio, <i>An</i>
32	[01]	introduction to statistical modeling of extreme
33		values, vol. 208. Springer, 2001.
34	[65]	V. Hodge and J. Austin, "A survey of outlier
35		detection methodologies," Artif. Intell. Rev., vol.
36		22, no. 2, pp. 85–126, 2004.
37	[66]	R. J. Hyndman, "The problem with Sturges' rule
38		for constructing histograms," <i>Monash Univ.</i> , no.
39	[67]	July, pp. 1–2, 1995. Ž Jugić A. I. Connolly, J. T. VanderPlac, and A.
40 11	[07]	L. IVEZIC, A. J. COIMONY, J. I. Vanuel Plas, and A. Cray Statistics Data Mining and Machine
41		Learning in Astronomy: A Practical Python Guide
43		for the Analysis of Survey Data. Princeton
44		University Press, 2014.
45	[68]	I. Salgado-Ugarte, M. Shimizu, and T. Taniuchi,
46		"Practical rules for bandwidth selection in
47		univariate density estimation," Stata Tech. Bull.,
48		vol. 5, no. 27, pp. 5–19, 1995.
49	[69]	D. Freedman and P. Diaconis, "On the histogram
50		as a density estimator:L2 theory," Zeitschrift für
51		Cabiata vol 57 no 4 nn 452 476 1091
52	[70]	D. I. C. Tinghui Wang "Towards Unsupervised
53	[/0]	Multi-Resident Tracking in Amhient Assisted
55		Living: Methods and Performance Metrics." in
56		Assistive Technology for the Elderly, 1st Edition, N.
57		S. Subhas Mukhopadhyay, Ed. Academic Press.
58	[71]	M. Rosenblatt, Markov Processes, Structure and
59		Asymptotic Behavior: Structure and Asymptotic
60		Behavior, vol. 184. Springer Science & Business

Media, 2012.

61

73

- K. Doty, S. Roy, and T. R. Fischer, "Filtering and smoothing state estimation for flag Hidden
 Markov Models," in *American Control Conference*(ACC), 2016, 2016, pp. 7042–7047.
- 66 [73] M. Xue and S. Roy, "Spectral and graph-theoretic
 67 bounds on steady-state-probability estimation
 68 performance for an ergodic Markov chain," J.
 69 Franklin Inst., vol. 348, no. 9, pp. 2448–2467,
 70 2011.
- 71 [74] S. Roy, "Scaled consensus," Automatica, vol. 51,
 72 pp. 259–262, 2015.

17