

Unsupervised Detection and Analysis of Changes in Everyday Physical Activity Data

Gina Sprint^{a,*}, Diane J. Cook^a, Maureen Schmitter-Edgecombe^b

^a*School of Electrical Engineering and Computer Science, Washington State University,
Pullman, WA*

^b*Department of Psychology, Washington State University, Pullman, WA*

Abstract

To be written later.

Keywords: Physical activity monitoring, Wearable sensors, Unsupervised learning, Change point detection, Data mining

1. Introduction

1 In recent years, sensors have become ubiquitous in our everyday lives.
2 Sensors are ambient in the environment, embedded in smartphones, and
3 worn on the body. Data collected from sensors form a time series, where
4 each sample of data is paired with an associated timestamp. This sensor-
5 based time series data is valuable when monitoring human behavior to detect
6 and analyze changes in behavior. Such analysis can be used to detect seasonal
7 variations, new family or job situations, or health events. Analyzing sensor-
8 based time series data can also be used to monitor changes in human behavior
9 as a person makes progress toward a fitness goal. Making a significant lifestyle
10 change often takes weeks or months of establishing new behavior patterns [1],
11 which can be challenging to sustain. Automatically detecting and tracking
12 behavior changes from sensor data can provide a valuable motivating and
13 monitoring tool.
14

15 Recently, wearable sensors have increased in popularity as people aspire
16 to be more conscientious of their physical health. Many consumers purchase

*Corresponding author

Email addresses: gsprint@eecs.wsu.edu (Gina Sprint), cook@eecs.wsu.edu
(Diane J. Cook), schmitter-e@wsu.edu (Maureen Schmitter-Edgecombe)

17 a pedometer or wearable fitness device in order to track their physical ac-
18 tivity, often in pursuit of a goal such as increasing cardiovascular strength,
19 losing weight, or improving overall health. Physical activity is estimated by
20 pedometers and fitness trackers in terms of the steps taken by the wearer
21 [2]. To track different types of changes in physical activity data, two or more
22 time periods, or windows, of physical activity data can be quantitatively and
23 objectively compared. If the two time windows contain significantly different
24 sensor data then this may indicate a significant behavior change. Existing
25 off-the-shelf change point detection methods are available to detect change
26 in time series data, but the methods do not provide context or *explana-*
27 *tion* regarding the detected change. For physical activity data, algorithmic
28 approaches to change detection require additional information about what
29 type of change is detected and its magnitude to potentially report progress to
30 users for motivation and encouragement purposes. Furthermore, existing ap-
31 proaches often do not provide a method for determining if a detected change
32 is *significant*, meaning the magnitude of change is high enough to suspect
33 it likely resulted from a lifestyle alteration. A personalized, data-driven ap-
34 proach to significance testing for fitness tracker users is a necessary feature
35 of physical activity change detection.

36 Currently, there is no clear consensus regarding which change detection
37 approaches are best for detecting and analyzing changes in physical activity
38 data. Consequently, we aim to formalize the problem of unsupervised physi-
39 cal activity change detection and analysis and address the problem with our
40 Behavior Change Detection (BCD) approach. BCD is a framework that 1)
41 segments time series data into time periods, 2) detects changes between time
42 periods 3) determines significance of the detected changes and 4) analyzes the
43 significant changes. We review recently proposed change detection methods
44 and analyze their appropriateness for BCD.

45 We demonstrate the approaches on sample FitBit data collected for two
46 weeks from an individual with documented daily activity information [TODO:
47 which data is this? not presently included]. Next, we evaluate the ability of
48 alternative change detection approaches to capture pattern changes in syn-
49 thetic physical activity data. Finally, we illustrate how change approaches
50 are used to monitor, quantify, and explain behavior differences in Fitbit data
51 collected from older adults who participate in a brain health behavior in-
52 tervention. We conclude with discussions about the limitations of current
53 approaches and suggestions for continued research on unsupervised sensor-
54 based change detection.

55 2. Related Work

56 In the literature, a few studies have aimed to detect change specifically in
57 human behavior patterns. These approaches have quantified change statisti-
58 cally [3, 4], graphically [4, 5, 6], and the change detection algorithms that
59 we incorporate into our method, BCD. [5, 7, 8, 9]. Recently, Merilahti *et*
60 *al.* [3] extracted features derived from actigraphy data collected for at least
61 one year. Each feature was individually correlated with a component of the
62 Resident Assessment Instrument for insights into how longitudinal changes
63 in actigraphy and functioning are associated. While this approach provides
64 insight into the relationship between wearable sensor data and clinical as-
65 sessment scores, this study does not directly quantify sensor-based change.

66 Wang *et al.* [5] introduced another activity-based change detection ap-
67 proach in which passive infrared motion sensors were installed in apartments
68 and utilized to estimate physical activity in the home and time away from
69 home. This data were converted into gray-level co-occurrence matrices for
70 computation of image-based texture features. Their case studies suggest the
71 proposed texture method can detect lifestyle changes, such as knee replace-
72 ment surgery and recovery. Though the approach does not provide expla-
73 nation of the detected changes over time, visual inspection of the data is
74 suggested with activity density maps. More recently, Tan *et al.* [6] applied
75 the texture method to data from Fitbit Flex sensors for tracking changes in
76 daily activity patterns for elderly participants. Another approach for activity
77 monitoring is the Permutation-based Change Detection in Activity Routine
78 (PCAR) algorithm [7]. PCAR researchers modeled activity distributions for
79 time windows containing at least one day. Changes between windows were
80 quantified with a probability of change acquired via hypothesis testing.

81 The change detection algorithms described previously are intended for
82 monitoring human activity behavior. There are several additional approaches
83 that are not specific to activity data, but instead represent generic statistical
84 approaches to detecting changes in time series data. Change point detection,
85 the problem of identifying abrupt changes in time series data [10], consti-
86 tutes an extensive body of research as there are many applications requiring
87 efficient, effective algorithms for reliably detecting variation. There are many
88 families of change detection algorithms that are suitable for different appli-
89 cations [11], for example: single point or two sample (window), univariate or
90 multivariate, labeled (supervised) or unlabeled (unsupervised), streaming or
91 offline, etc. Algorithms appropriately handling two sample, univariate, and

92 unlabeled data are most relevant to the current study due to their data-driven
93 change score computation and no need for ground truth information. Uni-
94 variate unsupervised change detection approaches include subspace models
95 and likelihood ratio methods [8]. One particular subgroup of likelihood ratio
96 methods, direct density ratio estimator methods, are used in various appli-
97 cations [12, 13]. Relative Unconstrained Least-Squares Importance Fitting
98 (RuLSIF) [8] is one such approach used to measure the difference between
99 two samples of data surrounding a candidate change point. Other recent
100 change point detection research includes work on multivariate [14, 15] and
101 streaming time series data [11].

102 The above approaches are effective methods for detecting change between
103 two samples of data; however, they are not explanatory methods as they only
104 identify if two samples are different and do not provide information on how
105 the samples are different. Once a change is detected and determined signif-
106 icant, additional analyses are required to explain the change that occurred.
107 Hido *et al.* [9] formalized this problem as change analysis, a method of anal-
108 ysis beyond change detection to explain the nature of discrepancy. Hido and
109 colleagues’ solution to change analysis utilizes supervised machine learning
110 algorithms to identify and describe changes in unsupervised data. Research
111 by Ng and Dash [16] and Yamada *et al.* [10] have also explored methods
112 for detecting and explaining change in time series data.

113 The aforementioned methods provide several options for change detection
114 and analysis, each with their own suitability for various applications (i.e. uni-
115 variate vs. multivariate, change significance testing available or not, etc.).
116 In this paper, the appropriateness of 1) RuLSIF [8], 2) texture-based dis-
117 similarity [5], 3) PCAR [7], 4) our proposed adaptation of PCAR for small
118 window sizes (sw-PCAR), and 5) change analysis [9] for unsupervised change
119 detection and analysis in wearable sensor time series data is evaluated.

120 3. Methods

121 Physical activity is often defined as any bodily movement by skeletal
122 muscles that results in caloric energy expenditure [17]. Physical activity
123 consists of bouts of movement that are separated by periods of rest. Physical
124 activity bouts are composed of four dimensions [17]:

- 125 1. Frequency: the number of bouts of physical activity within a time
126 period, such as a day.

- 127 2. Duration: the length of time an individual participates in a single bout.
 128 3. Intensity: the physiological effort associated with a particular type of
 129 physical activity bout.
 130 4. Activity type: the kind of exercise performed during the bout.

131 To add exercise throughout the day, individuals can increase their number
 132 of bouts (frequency), increase the length of bouts (duration), increase the
 133 intensity of bouts, and vary the type of physical activity performed during
 134 the bouts. These four components of physical activity represent four distinct
 135 types of changes that can reflect progress towards many different health goals,
 136 such as increasing physical activity or consistency in one’s daily routine.

137 We study the problem of detecting and analyzing change in physical activ-
 138 ity patterns. More specifically, we introduce methods to determine if a
 139 significant change exists between two windows of time series step data. Al-
 140 gorithm 1, BehaviorChangeDetection, outlines this process. Let m denote
 141 the number of equal sized time intervals in a day and t_{mins} denote the number
 142 of minutes per time interval. For example, if the sampling rate of the wear-
 143 able sensor device is one reading per minute, $t_{mins} = 1$ minute and $m = 1440$
 144 minutes / t_{mins} . Now, let $D = \{x_1, x_2, \dots, x_t, \dots, x_m\}$ be one day of time series
 145 data where x_t is a scalar number of steps taken at time interval $t = 1, 2, \dots, m$.
 146 Let $W_i = \{D_i, D_{i+1}, \dots, D_n\}$ be a window of n days with $1 \leq i \leq n$. Suppose
 147 we have two windows of data, W_i and W_j ($i \leq j$), sampled from a time se-
 148 ries X . For change detection and analysis, a function F computes a change
 149 score, $CS = F(W_i, W_j)$ between two windows, W_i and W_j . Furthermore, an
 150 aggregate window, \hat{W} , represents the average of all days within the window
 151 W :

$$\hat{W} = \frac{1}{n} \sum_{i=1}^n D_i, D_i \in W \quad (1)$$

152 We can compare windows of data within time series data X . These win-
 153 dows may represent consecutive times (e.g., days, weeks, months), a baseline
 154 window (e.g., the first week) with each subsequent time window, or overlap-
 155 ping windows. Windows W_i and W_j can be formed as subsets of X based
 156 on an *offset* denoting the start of W_j as a function of the start of W_i and
 157 iteration advancements adv_i and adv_j to move windows W_i and W_j respec-
 158 tively for the next comparison. Two windows can be compared in either
 159 baseline or sliding window mode. For a baseline window comparison, the
 160 first window is a reference window that occurs at the beginning of the time
 161 series (i is initialized to 1) and is used in each comparison, so $adj_i = 0$. All

162 subsequent windows are compared to the baseline window. Thus j is initial-
 163 ized to $1 + offset$ and is subsequently advanced by adv_j . In the case of a
 164 sliding window comparison, both windows used for comparison are advanced
 165 through the time series data. Typically $adv_i = adv_j$ for consistently spaced
 166 comparisons. In Algorithm 1, BehaviorChangeDetection, i is initialized to
 167 0 and j is initialized to $offset$. In steps 17 and 18, i is advanced to $i + adv_i$
 168 and j is advanced to $j + adv_j$.

Algorithm 1 BehaviorChangeDetection($X, n, offset, adv_i, adv_j$)

```

1: Input:  $X$  = time series data
2: Input:  $n$  = window length in days
3: Input:  $offset$  = number of days separating windows
4: Input:  $adv_i$  = number of days to advance the first window
5: Input:  $adv_j$  = number of days to advance the second window
6: Output: Change score vector  $V$ 
7: Initialize:  $i = 1$  and  $j = 1 + offset$ 
8: for each pair of windows to compare,  $W_i$  and  $W_j$  of time series  $X$ :
9:    $W_i = X[i : i + n]$ 
10:   $W_j = X[j : i + n]$ 
11:  Compute  $CS = F(W_i, W_j)$ 
12:  Determine if  $CS$  is significant
13:  Identify the type of change that is exhibited
14:    Manual inspection of change
15:    Unsupervised inspection (change analysis)
16:  Append  $CS$  to change score vector  $V$ 
17:   $i = i + adv_i$ 
18:   $j = j + adv_j$ 
end for
19: return Change score vector  $V$ 

```

169 The choice of window size, n , limits the algorithms that can be applied to
 170 the data. For example, the PCAR algorithm [7] is designed for longitudinal
 171 data comprising several months; consequently sensitivity decreases with small
 172 window sizes. For this study, we categorize window size n choices into the
 173 following descriptors:

- 174 1. Small window ($n = 1$ day). Suitable for performing day-to-day com-
 175 parisons (ex: $D_{Monday1}$ compared to $D_{Monday2}$, $D_{Tuesday1}$ compared to

- 176 $D_{Tuesday2}, \dots$) or aggregate day comparisons (ex: \hat{W}_1 compared to \hat{W}_2 ,
177 \hat{W}_2 compared to \hat{W}_3, \dots).
- 178 2. Medium window ($2 \text{ days} \leq n \leq 5 \text{ days}$). Suitable for performing
179 weekday-to-weekday (ex: W_1 compared to W_2 where $W_1 = \{D_{Monday1},$
180 $D_{Tuesday1}, D_{Wednesday1}, D_{Thursday1}, D_{Friday1}\}$ and $W_2 = \{D_{Monday2},$
181 $D_{Tuesday2}, D_{Wednesday2}, D_{Thursday2}, D_{Friday2}\}$ or weekend-to-weekend
182 comparisons.
- 183 3. Large window ($n > 5 \text{ days}$). Suitable for performing week-to-week or
184 month-to-month comparisons.

185 3.1. Change Detection Algorithms

186 In the following sections, we describe five different algorithmic options for
187 the window-based change score function, F . A summary and comparison of
188 the algorithms is listed in Table 1.

189 3.1.1. RuLSIF

190 Non-parametric approaches to change point detection include a family of
191 methods comparing the probability distributions of two time series samples
192 to determine the corresponding dissimilarity. A greater difference between
193 the two distributions implies a higher likelihood that change occurred be-
194 tween the two samples. Instead of estimating the probability distributions,
195 their ratio can be estimated and used to detect changes in the underlying
196 probability distributions. Direct density ratio estimation between two win-
197 dows of time series data is substantially simpler to solve than computing the
198 windows' probability densities independently and then using these to com-
199 pute the ratio. Unconstrained Least-Squares Importance Fitting (uLSIF) [8]
200 is one such ratio estimation approach that measures the difference between
201 two samples of data surrounding a candidate change point. For this ap-
202 proach, the density ratio between two probability distributions is estimated
203 directly with the Pearson divergence dissimilarity measure. Depending upon
204 the data, the Pearson divergence can be unbounded. Consequently, a modifi-
205 cation to uLSIF, relative uLSIF (RuLSIF), utilizes an alpha-relative Pearson
206 divergence to bound the change score by $1/\alpha$, $0 \leq \alpha < 1$ [8].

207 3.1.2. Texture-based Change Detection

208 For this approach, two windows of physical activity data, W_i and W_j ,
209 are converted into gray-level co-occurrence matrices (GLCM) [5]. Rows in
210 the resulting GLCM correspond to time intervals while columns correspond

Table 1: Window-based change detection algorithms.

Approach	Window size	Preprocessing	Change score	Change significance test
RuLSIF [8]	Any	Optional Hankel matrix [8]	Probability density ratio estimation with Pearson divergence	Threshold learning in supervised applications. N/A for unsupervised applications
Texture-based [5, 6]	Any	Gray-level co-occurrence matrix, texture features	Weighted normalized Euclidean distance	N/A
PCAR [7]	Large	$m \times N$ KL divergence permutation matrices	Count of time intervals with significant changes (proportion of permuted KL distances greater than observed window)	N/A
sw-PCAR	Small, medium	N KL divergence permutation vectors	KL divergence distance	Non-parametric outlier detection based on Boxplot analysis
Virtual classifier [9]	Large	Physical activity features (intraday and interday if window size > 1)	Cross validation prediction accuracy of binary classifier	Hypothesis testing based on prediction accuracy exceeding a threshold

N = number of permutations

211 to days. Each cell of the GLCM contains normalized step values symbol-
 212 ized by a gray scale value [TODO: possibly add example (see Figure 4 from
 213 Tan)]. Next, the texture features of contrast, dissimilarity, homogeneity, an-
 214 gular second moment (ASM) energy, and correlation are computed from each
 215 GLCM [18], producing feature vectors T_i and T_j . To compare two windows
 216 W_i and W_j for changes, a weighted normalized Euclidean distance measure is
 217 used as a change score to quantify the differences between the corresponding
 218 feature vectors T_i and T_j . The smaller the Euclidean distance between these
 219 two vectors, the more similar the two windows of data are. The texture-based
 220 approach can operate on small or large window sizes; however, the method
 221 lends itself more appropriately to large window sizes (Wang *et al.* [5] used
 222 window size of one month).

223 3.1.3. PCAR

224 PCAR utilizes smart home sensor data to detect changes in behavioral
 225 routines [7]. This approach assumes that an activity recognition algorithm
 226 [19] is available to label the sensor data with corresponding activity names A
 227 $= \{A_1, A_2, \dots, A_a\}$. The algorithm is based on the notion of an activity curve
 228 C , a compilation of m probability distributions $R_t = \{d_{t,1}, d_{t,2}, \dots, d_{t,a}\}$ of a
 229 activities per time interval t in a day ($t = 1, 2, \dots, l, \dots, m$). The l^{th} element
 230 $d_{t,l}$ represents the probability of activity A_l during time interval t . Windows
 231 of time spanning multiple days are averaged into an aggregate activity curve
 232 \hat{A} . To compute a change score CS between two aggregated activity curves
 233 \hat{A}_i corresponding to window W_i and \hat{A}_j corresponding to window W_j , the
 234 two activity curves are first maximally aligned with dynamic time warping
 235 (DTW). Next, the symmetric Kullback-Leibler (KL) divergence is used to
 236 compute the distance between each pair of activity distributions in A_i and
 237 A_j at time interval t [7]:

$$KL_{symmetric} = KL(D_i, D_j) + KL(D_j, D_i) \quad (2)$$

238 where

$$KL(D_i, D_j) = \sum_{k=1}^a d_{i,k} \cdot \log \frac{d_{i,k}}{d_{j,k}} \quad (3)$$

239 The total distance between the two curves is calculated as the sum of
 240 each time interval distance. To test significance of the activity curve distance
 241 value, W_i and W_j are concatenated to form a window W of length $2n$ days.

242 Next, all days within W are shuffled. The first half of the shuffled days form a
 243 new first window, W_i^* , while the second half form a new second window, W_j^* .
 244 KL distances for each DTW-aligned time interval pairs in W_i^* and W_j^* form a
 245 vector that is inserted into a matrix. This shuffling procedure is repeated N
 246 times, producing a $N \times m$ permutation matrix, M . If N is large enough, M
 247 forms an empirical distribution of the possible permutations of activity data
 248 within the two windows of time. Next, for each time interval t , the number of
 249 permuted KL distances that exceed the original change score CS is divided
 250 by N to form a p-value. After computing a p-value for each time interval t ,
 251 the Benjamini-Hochberg correction [20] is applied for a given α ($\alpha = 0.01$
 252 or 0.05). Finally, the remaining significant p-values are counted to produce
 253 the change score. For this study, we additionally normalize the change score.
 254 To do this, we divide the score by the result of DTW-alignment pairs and
 255 multiplying by 100 to yield a percentage value.

256 While the algorithm is intended for activity distribution data available
 257 from activity recognition algorithms, in this paper we adapt the PCAR algo-
 258 rithm to analyze the physical activity change detection as part of our BCD
 259 method. Instead of activity distribution vectors, we use scalar step counts.
 260 Additionally, PCAR is suitable for only large window sizes due to the re-
 261 quirement of permuting daily time series data. The approach originally was
 262 intended for correlating change scores with standardized clinical assessments
 263 to determine if ambient smart home sensor-based algorithms can detect cog-
 264 nitive decline [7]. Consequently, there is not a test for significance of the
 265 change score. In the following section we propose a version of PCAR that
 266 is more suitable for small windows (sw-PCAR) as required by BCD and in
 267 Section 3.2 we propose an accompanying significance test for sw-PCAR.

268 3.1.4. *sw-PCAR*

269 We propose a small window adaptation of PCAR to allow permutation-
 270 based change detection for window sizes of one week or less. Algorithm
 271 2 outlines the sw-PCAR approach. For sw-PCAR, two windows W_i and
 272 W_j are collapsed into aggregate windows \hat{W}_i and \hat{W}_j (see Equation 1). A
 273 change score CS is derived by computing the KL divergence between the
 274 average number of steps taken in \hat{W}_i and the average number of steps taken
 275 in \hat{W}_j . Next, \hat{W}_i and \hat{W}_j are concatenated to form a window W of length
 276 two days. All time intervals within W are shuffled. The first half of the
 277 shuffled intervals form a new first window, W_i^* , while the second half form a
 278 new second window, W_j^* . W_i^* and W_j^* are each averaged to produce two step

279 values. The KL distance between the two values is computed and inserted
 280 into a vector. This is repeated N times to produce a N -length vector V of
 281 KL distances. Vector V is later used for change score significance testing (see
 282 Section 3.2).

Algorithm 2 sw-PCAR(W_i, W_j, N)

1: Input: W_i, W_j = two windows of time series data
 2: Input: N = number of permutations
 3: Output: Change score CS and Boolean sig
 4: Initialize: $k = 0$
 5: Initialize: V as a vector of length N
 6: Compute \hat{W}_i, \hat{W}_j aggregate windows
 7: Compute CS , the KL divergence between \hat{W}_i and \hat{W}_j
 8: **while** $k < N$:
 9: Shuffle the time intervals of \hat{W}_i and \hat{W}_j
 10: Generate new aggregate windows W_i^* and W_j^*
 11: Compute KL divergence between W_i^* and W_j^*
 12: Store resulting distance in V
 13: $k = k + 1$
end while
 14: $sig = \text{BoxplotOutlierDetection}(CS, V)$ (see Algorithm 3)
 15: **return** CS, sig

283 *3.1.5. Virtual Classifier*

284 Change analysis, as proposed by Hido *et al.* [9], utilizes a virtual binary
 285 classifier to detect and investigate change. We apply the virtual classifier
 286 (VC) approach to the physical activity change problem for large window sizes.
 287 First, a feature extraction step reduces two windows W_i and W_j into two $n \times z$
 288 feature matrices, M_i and M_j , where n is the window size (in days) and z is the
 289 number of features that are extracted. Next, each daily feature vector of M_i
 290 is labeled with a positive class and each daily feature vector of M_j is labeled
 291 with a negative class. VC trains a decision tree to learn the decision boundary
 292 between the virtual positive and negative classes. The resulting average
 293 prediction accuracy based on k -fold cross validation represented as p_{VC} . If
 294 a significant change exists between W_i and W_j , the average classification
 295 accuracy p_{VC} of the learner should be significantly higher than the accuracy
 296 expected from random noise, $p_{rand} = 0.5$, the binomial maximum likelihood
 297 of two equal length windows.

298 *3.2. Change Significance Testing*

299 Significance testing of change score CS is necessary to interpret change
300 score values. For the VC approach, Hido *et al.* [9] proposed a test of signif-
301 icance to determine if p_{VC} is significantly greater than p_{rand} . For this test,
302 the inverse survival function of a binomial distribution is used to determine
303 a critical value, $p_{critical}$, at which n Bernoulli trials are expected to exceed
304 p_{bin} at α significance. If $p_{vc} \geq p_{critical}$, a significant change exists between
305 the two windows, W_i and W_j .

306 The PCAR approach does not have an accompanying test of significance.
307 We address this with our proposed sw-PCAR technique. sw-PCAR com-
308 puts change significance by comparing CS to the permutation vector V
309 with boxplot-based outlier detection (see Algorithm 3). An outlier can be
310 defined as an observation which appears to be inconsistent with other obser-
311 vations in the dataset [21]. For this method, the interquartile range (75th
312 percentile - 25th percentile) of V is computed. Values outside of the $1.5 \cdot$
313 75th percentile are considered outliers [22]. If CS is determined to be an
314 outlier of V , then the change score is considered significant. There are al-
315 ternative approaches to test membership of an observation (i.e. CS) to a
316 sample distribution (i.e. V) other than boxplot outlier detection. If the sam-
317 ple is normal, statistical tests such as Grubb’s test for outliers [23] can be
318 applied. However, the assumption of normality does not hold for all samples
319 of human behavior data. More advanced alternatives include data mining
320 techniques relevant to outlier detection [21, 24]. Exploration and testing of
321 such data mining techniques are outside the scope of this paper.

322 RuLSIF does not explicitly provide a method to determine a cutoff thresh-
323 old for values of the Pearson divergence function are considered significant
324 change scores. In supervised applications where ground truth change labels
325 are available, a threshold parameter is typically learned by repeated training
326 and testing with different parameter values. For unsupervised applications,
327 domain knowledge and/or alternative data-driven approaches are necessary.
328 Like RuLSIF, the texture-based method also does not provide a test of change
329 significance. For RuLSIF and texture-based approaches, we propose a medium
330 to large window change significance test based on intra-window variability
331 and outlier detection.

332 Our proposed change significance test utilizes the existence of day-to-day
333 variability in human activity patterns [25]. In order to consider a change
334 between two windows significant, the magnitude of change should exceed
335 the day-to-day variability within each window. To illustrate, consider two

Algorithm 3 BoxplotOutlierDetection(CS, V)

```
1: Input:  $CS$  = change score between two windows
2: Input:  $V$  = sample distribution vector
3: Output: Boolean  $sig$ 
4: Arrange  $V$  in ascending order
5: Compute  $Q_1$ , the 25th percentile of  $V$ 
6: Compute  $Q_3$ , the 75th percentile of  $V$ 
7: Compute the interquartile range of  $V$ ,  $IQR = Q_3 - Q_1$ 
8: if  $CS > 1.5 \cdot IQR$ :
9:    $sig = True$ 
10: else:
11:    $sig = False$ 
12: return  $sig$ 
```

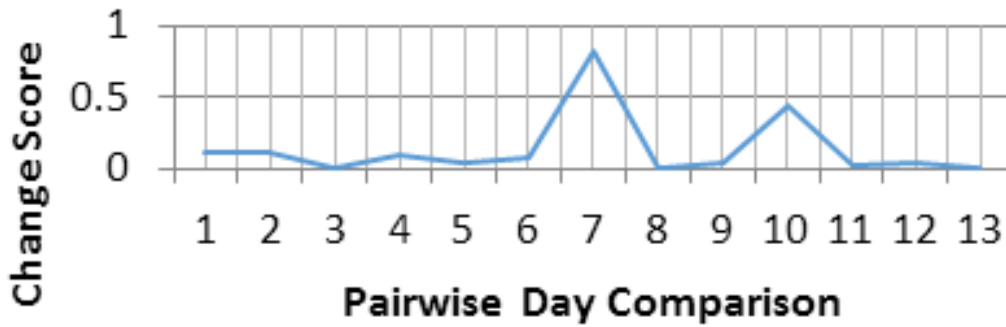


Figure 1: Pairwise sliding window RulSIF change scores. [TODO: MAKE THIS A MATPLOTLIB PDF]

336 adjacent, non-overlapping windows W_1 and W_8 , each of length $n = 7$ days.
337 Now run a pairwise sliding window change algorithm over W_1 concatenated
338 with W_8 . If there is a significant change between the windows, the magnitude
339 of change should be higher for the inter-window comparison (between days 7
340 and 8) than any other intra-window comparison. Fig. 1 shows an example
341 plot of RulSIF change scores for real Fitbit data illustrating this phenomenon.
342 There are small, noisy day-to-day changes for all comparisons except the
343 largest maximum occurring for inter-window comparison (7th change score)
344 and a potential anomaly between days 3 and 4 of W_8 .

345 Based on the assumption that a significant inter-window change should
346 exceed intra-window change, we propose an intra-window change significance

347 test (see Algorithm 4). Given a change score CS between two windows, the
 348 task is to determine if CS is significant. To do this, first compute a list of
 349 all possible daily change scores, DCS , within each window. DCS contains
 350 $2 \cdot \text{Combination}(n, 2)$ change scores (see Algorithm 5). For example, a
 351 week to week comparison ($n = 7$) would generate an intra-window daily
 352 change score sample of 42 day-to-day variations. Next, apply the outlier
 353 detection method (see Algorithm 3) from sw-PCAR to test if CS is an outlier
 354 score when compared to the distribution of intra-window daily change scores
 355 DCS . Advantages of the proposed test include it is non-parametric and
 356 can be coupled with any small window change algorithms. Furthermore, the
 357 approach lends itself well to online change detection algorithms, since only a
 358 vector of the previous window's (baseline or sliding) change scores need to be
 359 retained. Finally, the candidate change score, CS , can be computed based
 360 on any window size (i.e. Monday to Monday, aggregate to aggregate, week
 361 to week, etc.).

Algorithm 4 Intra-windowChangeSignificanceTest(W_1, W_2, n, CS, F)

- 1: Input: $W_1, W_2 =$ two windows of time series data
 - 2: Input: $n =$ window size
 - 3: Input: $CS =$ change score between W_1 and W_2
 - 4: Input: $F =$ change score function
 - 5: Output: Boolean sig
 - 6: Initialize: Vector of daily change scores DCS
 - 7: Append W_1 intra-window daily change scores to DCS (see Algorithm 5)
 - 8: Append W_2 intra-window daily change scores to DCS (see Algorithm 5)
 - 9: Compute $sig = \text{BoxplotOutlierDetection}(CS, DCS)$ (see Algorithm 3)
 - 10: **return** sig
-

362 *3.3. Change Analysis*

363 If a change significance test concludes a change score is significant, the
 364 next step is to determine the source of change (see Algorithm 1 for an
 365 overview of the change detection and analysis process). Often this step re-
 366 quires the computation of features that summarize the data and provide a
 367 meaningful context for change. For example, the number of daily steps taken
 368 is an example of a simple physical activity feature. The change between daily
 369 steps from one window of time to the next can be quantified and used for
 370 an explanation of change. Several approaches exist to capture change across

Algorithm 5 Intra-windowDailyChangeScores(W, n, F)

```
1: Input:  $W$  = window of time series data
2: Input:  $n$  = window size
3: Input:  $F$  = change score function
4: Output: Vector of daily change scores  $DCS$ 
5: Initialize:  $i = 0, j = 0$ 
6: while  $i < n - 1$ :
7:    $curr = W[i]$ 
8:    $j = i + 1$ 
9:   while  $j < n$ :
10:     $next = W[j]$ 
11:     $CS = F(curr, next)$ 
12:    Append  $CS$  to  $DCS$ 
13:     $j = j + 1$ 
14:   end while
15:    $i = i + 1$ 
16: end while
17: return  $DCS$ 
```

371 time in individual metrics. A straightforward method is to compute the per-
372 cent change for a feature f from a previous window W_1 to a current window
373 W_2 :

$$\Delta\% = \frac{f_{W_2} - f_{W_1}}{f_{W_1}} \quad (4)$$

374 Statistical approaches such as two sample tests or effect size analyses can
375 also be applied to quantify change; however, in applying repeated statistical
376 tests, the multiple testing problem should be accounted for with a method
377 such as the Bonferroni or Benjamini-Hochberg correction [20].

378 One of the advantages of the virtual classifier approach over other change
379 point detection algorithms is it includes an explanation of the source of
380 change without reliance on statistical tests. Upon significant change detec-
381 tion, retraining a decision tree on the entire dataset and inspecting the tree
382 reveals which features are most discriminatory in learning the differences be-
383 tween two windows. Naturally, this approach requires a pre-processing step
384 to compute relevant features from the windowed time series data. The fol-
385 lowing section presents relevant features utilized for physical activity change
386 analysis.

387 *3.4. Feature Extraction*

388 The following features are grouped together based on the number of days
389 required for computation: 1) one day (24 hour window of time) or less, 2)
390 at least one day, or 3) two or more days. Daily features include intraday
391 physical activity summaries based on intensity, frequency, duration, variabil-
392 ity of steps and walking bouts. Sequences of time series data with steps
393 greater than t_{mins} . t_{mins} represents the minimum number of steps per time
394 interval to be considered physical activity. This assumes physical activity
395 is characterized by at least one step per minute. If ground truth activity
396 labels, such as walking, biking, chores, etc., are available from the device
397 user and/or an activity recognition algorithm, physical activity type can be
398 inferred and thresh can be set dynamically for different activities. For this
399 study, we assume such labeled information is not available and set thresh to
400 t_{mins} .

- 401 • Daily PA intensity
 - 402 – Bout steps: Mean and SD of number of steps per bout
 - 403 – Period steps [4]: Mean and SD per period 1) 24 hour period (full
 - 404 days), 2) Day (9am-9pm), 3) Night (12am-6am). Day and night
 - 405 normalized by 24 hour mean
 - 406 – Ratio of mean night and day steps [4]: See period steps definition
- 407 • Daily PA frequency
 - 408 – Number of bouts: Count of detected PA bouts
- 409 • Daily PA duration
 - 410 – Bout minutes: Mean and SD of duration of bouts
 - 411 – PA intensity percentage: Mean percentage of 1) sedentary (< 5
 - 412 steps/min), 2) low ($5 \leq$ steps/min < 40), 3) moderate ($40 \leq$
 - 413 steps/min < 100), 4) high (\geq 100 steps/min) activity levels
 - 414 – Rest minutes: Mean and SD of duration of rest periods

415 Features computed on window sizes of at least one day include an adap-
416 tation of relative amplitude from Merilahti *et al.* [3] and texture features
417 from the texture-based change detection approach [18] (see Section 3.1.2).

418 • Relative amplitude: Normalized ratio between the most active 8 hours
 419 and the least active 4 hours activity periods (not required to be consec-
 420 utive). If sleep data is available, awake hours are used for least active
 421 periods. $RA = \frac{M8-L4}{M8+L4}$

422 • Texture features: See section 3.1.2

423 Features requiring at least two days of data summarize activity across or
 424 between days or quantify the users circadian rhythm (the periodicity from
 425 day-to-day [25]). Poincare-plot analysis [4] provides an additional set of
 426 useful physical activity features. Poincare plots depict how activity patterns
 427 repeat themselves based on a time delay, d . Time series data at time t ,
 428 $A(t)$, is plotted as a function of previous data, $A(t - d)$. From the resulting
 429 Poincare plot, two measures of dispersion can be computed 1) SD1, the
 430 standard deviation of the data against the axis $x = y$ and 2) SD2, standard
 431 deviation of the data against the axis orthogonal to $x=y$ and crosses this axis
 432 at the mean value of the data (center of mass). Delay values of $d = 24hours$
 433 and $d = 12hours$ plot the data as a function of the previous data and the
 434 counter phase respectively. The day-to-day circadian rhythm preservation
 435 (CRP) feature is based on dispersion values from these two delays. [TODO:
 436 include sample poincare plots?]

437 • Inter-daily stability (IS) [3]: Quantifies stability between the days
 438 $IS = \frac{n\sigma_{h=1,p}(x_h-x)^2}{p \sum_{i=1}^n (x_i-x)^2}$

439 • Intra-daily variability (IV) [3]: Quantifies the fragmentation of rhythm
 440 and activity $IV = \frac{n \sum_{i=2}^n (x_i-x_{i-1})^2}{(n-1) \sum_{i=1}^n (x_i-x)^2}$

441 • Circadian rhythm strength (CRS) [3]: Divides average night-time ac-
 442 tivity (11pm- 5am) by the average activity of the previous day (8am-
 443 8pm) $CRS = \frac{steps_{11pm-5am}}{steps_{prev8am-8pm}}$

444 • Cosinor mesor: Time series mean from fitting a cosinor functional
 445 model with a 24 hour period to time series data via least squares
 446 method [3, 4]

447 • Cosinor amplitude: Difference between the mesor and peak (or trough)
 448 of the fitted waveform

449 • Cosinor acrophase: Time of day at which the peak of a rhythm occurs

- 450 • Poincare plot SD1 [4]: Standard deviation of Poincare data against
451 the axis $x = y$
- 452 • Poincare plot SD2 [4]: Standard deviation of Poincare data against
453 the axis orthogonal to $x = y$ and crosses this axis at the mean value of
454 the data (center of mass)
- 455 • Poincare plot circadian rhythm preservation (CRP) [4]: Day-to-day
456 circadian rhythm preservation based on dispersion values from SD1 and
457 SD2 with delays of 24 hours and 12 hours $CRP = SD2_{24h} + SD1_{12h} -$
458 $SD1_{24h} - SD2_{12h}$

459 3.5. Datasets

460 To evaluate the change detection algorithms, two datasets are presented,
461 Hybrid-synthetic (HS) and B-Fit (BF). The HS dataset comprises synthetic
462 data and the BF dataset comprises real-world data collected from a Fitbit
463 study. To generate the HS dataset, step data collected from a volunteer wear-
464 ing a Fitbit Charge Heart Rate fitness tracker was re-sampled and modified
465 to produce five different synthetic physical activity profiles, each exhibiting
466 a different type of change. The length of HS profiles was set to 12 days,
467 resulting in two equal size windows of 6 days for comparison. Twelve days
468 was chosen for similarity to the real-world BF dataset. The HS profiles with
469 their profile identification (HS0-4) and a description are as follows:

- 470 1. HS0: No significant daily or window change. Data is subject to small
471 daily variation.
- 472 2. HS1: Medium daily change and consequently significant window change.
473 Increase bout duration and intensity from day-to-day.
- 474 3. HS2: No significant daily change but significant window change. In-
475 creased activity for days 7-12.
- 476 4. HS3: Medium daily change and consequently significant window change.
477 Increase activity variability from day-to-day.
- 478 5. HS4: No significant daily change for days 1-6. Significant daily change
479 for days 7-12. Consequently significant window change.

480 Figure 2 shows the associated activity density maps (ADMs) for HS1-4
481 profiles. An ADM is a heat map proposed by Wang *et al.* [5] to visualize
482 daily activity (shade of color) as a function of 24 hour time (Y-axis) and
483 window time (X-axis).

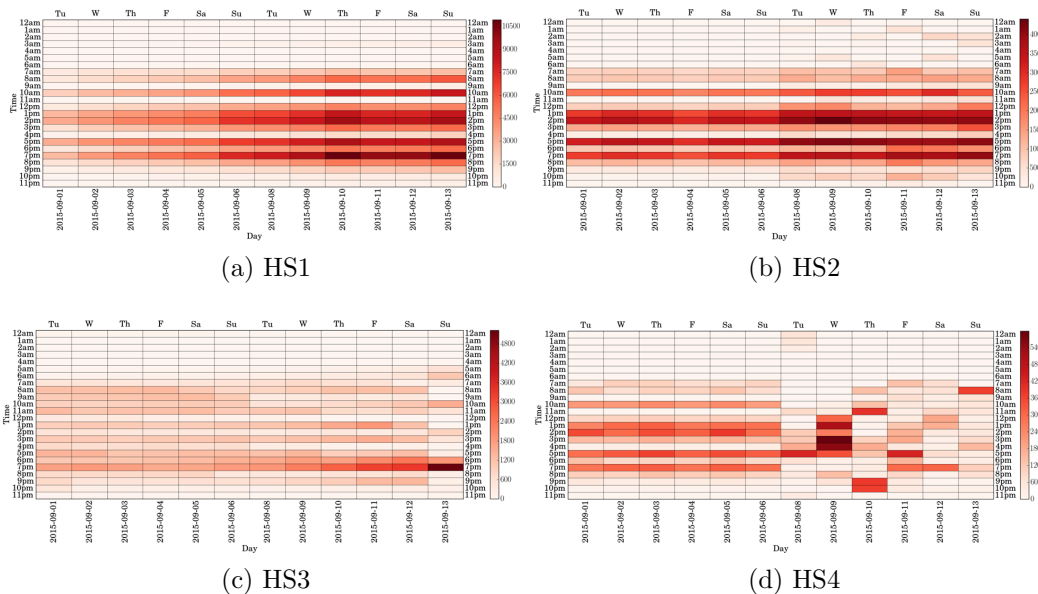


Figure 2: Hybrid-synthetic (HS) activity density maps for HS1-4.

484 The BF dataset consists of data collected from 11 older adults ($57.09 \pm$
485 8.79) participating in a 10-week health intervention called B-Fit (see Table 2
486 for participant characteristics). For this study, participants pre-intervention
487 physical activity profiles were assessed with wrist-worn Fitbit Flex fitness
488 trackers for one week (six full 24 hour days) before the intervention. During
489 weeks two through nine, the participants were educated in eight different
490 subjects related to health: exercise, cardiovascular risk factors, sleep, stress
491 reduction, cognitive engagement, nutrition, social engagement, and compen-
492 satory strategy use. Each week was devoted to education, goal setting, and
493 goal tracking for each of the eight categories. Example goals the partici-
494 pants set included “brisk walking, four times a week for a half hour or more”
495 and “drinking more water at work.” To track goal achievement, individuals
496 rated themselves from 0 to 3 (0: did not meet goal, 1: partly met goal, 2:
497 completely met goal, 3: exceeded goal). Table 2 shows self-ratings for the
498 BF participants for the categories of exercise and cardiovascular risk factors.
499 After the intervention, participants’ post-intervention physical activity pro-
500 files were assessed for one week (six full 24 hour days) with the same Fitbit
501 devices.

Table 2: B-Fit study participant characteristics

ID	Sex	Age	Height	Weight	S1 dates	S2 dates	Exercise	Cardio
BF1	F	65	5'4"	210 lbs	9/13-9/20	11/15-11/22	Brisk walking	2.5
BF3	F	68	5'3"	207 lbs	9/13-9/20	11/15-11/22	Walk more	2.86
BF18	F	65	5'2"	165 lbs	9/23-9/30	12/2-12/9	Take stairs	2
BF20	M	64	6'0"	190 lbs	9/23-9/30	12/2-12/9	Walking	2
BF23	F	57	5'5"	135 lbs	9/24-9/30	12/2-12/9	Yoga	1
BF24	F	53	5'6"	130 lbs	9/23-9/30	12/2-12/9	Walking	2
BF26	F	65	5'8"	145 lbs	9/23-9/30	12/2-12/9	Crossfit class	0
BF27	M	42	5'8"	145 lbs	9/23-9/30	12/2-12/9	Walking	3
BF28	M	48	6'4"	380 lbs	9/23-9/30	12/2-12/9	No exercise	3
BF29	F	52	5'6"	180 lbs	9/23-9/30	12/2-12/9	Walking	1
BF30	F	49	5'6"	262 lbs	9/23-9/30	12/2-12/9	Yoga	2

BF = B-Fit, ID = identification, lbs = pounds.

502 HS and BF data were subject to cleaning prior to serving as input to the
503 change detection algorithms. First, missing data were identified and handled
504 for BF data. Since donned/doffed information is not available for Fitbit Flex
505 fitness trackers, days with zero steps taken during the day (9am-9pm) were
506 considered missing data. By this method, three days exhibited missing data,
507 each belonging to a different participant (11/18/15: BF3; 9/28/15: BF27;
508 12/7/15: BF29, see Table 2 for participant information). Zero steps during
509 the day is likely due to removing the Fitbit to charge it, then forgetting to
510 put it back on until much later.

511 Algorithm 6, WindowedFillMissingData, summarizes the missing data
512 filling process. First, to fill a missing day ($D_{missing}$), the same day of the week
513 (i.e Monday) of $D_{missing}$ is identified in the opposite window (D_{other}). As an
514 example, P27 exhibited missing data on 9/28/15, which was a Monday in the
515 pre-intervention window. Next, the Monday of the post-intervention window
516 is identified (12/7/15 in this case). Euclidean distance-based clustering is
517 then applied to find the k nearest neighbors NN_{other} of D_{other} , $3 \leq k \leq 5$.
518 The days of the week for each day in NN_{other} are then identified. These are
519 used to select days, $NN_{missing}$ in the original window containing $D_{missing}$.
520 The k days of $NN_{missing}$ are averaged and used to fill $D_{missing}$ from 9am to
521 9pm. Three nearest neighbors is chosen as a minimum value for k in case
522 one of the days of the week is not available in the original six day window.

523 Additional pre-processing of Fitbit data includes down sampling the data
524 for a given time interval length, t_{mins} , by summing the steps every t_{mins} min-
525 utes. Furthermore, for the case of PCAR and sw-PCAR, add one smoothing
526 was applied to avoid a division by zero during KL divergence computations.
527 Finally, while the Fitbit Flex also provides distance traveled and calories
528 burned, for this study we only consider steps due to the high inter-metric
529 redundancy between steps and these two Fitbit metrics. For the BH data,
530 Pearson correlations of $r > 0.99$ for distance and $r > 0.90$ for calories were
531 measured.

532 4. Results

533 All data are processed with the Python programming language. Unless
534 otherwise stated, the following parameter values are used: k (fill missing
535 data): 3; n (window size): 6 days; $offset$ (window offset): 6 days; α (RuLSIF):
536 0.1; RuLSIF cross validation folds: 5; GLCM distance: 2; N (number of
537 PCAR and sw-PCAR permutations): 1000; α (PCAR): 0.05; virtual classifier

Algorithm 6 WindowedFillMissingData($D_{missing}, W_{missing}, W_{other}, k$)

- 1: Input: $D_{missing}$ = day with missing data to fill
 - 2: Input: $W_{missing}$ = window of time series data containing $D_{missing}$
 - 3: Input: W_{other} = window of time series data
 - 4: Input: k = number of nearest neighbors
 - 5: Output: D_{fill} day data to fill $D_{missing}$
 - 6: D_{other} = day in W_{other} with same day of week as $D_{missing}$
 - 7: NN_{other} = k nearest neighbors of D_{other} in W_{other}
 - 8: $NN_{missing}$ = k days in $W_{missing}$ with same days of week as days in NN_{other}
 - 9: $D_{missing}[9am:9pm]$ = average of k days in $NN_{missing}[9am:9pm]$
 - 10: **return** D_{fill}
-

538 cross validation folds: 4; virtual classifier prediction threshold $p_{critical}$: 0.75.
539 The time interval aggregation size t_{mins} is tested with values of $t_{mins} = 1, 5,$
540 $10, 15, \dots, 60$ minutes.

541 Table 3 shows RuLSIF, Texture-based, sw-PCAR, and VC significant
542 change results for each HS profile for each time interval length t_{mins} . Ta-
543 ble 4 shows PCAR change scores for each HS profile and BF participant
544 for each time interval length. The contextual features of number of bouts,
545 bout minutes, bout steps, daily steps, and sedentary minutes percent are
546 listed in Table 5 with window one and window two values (mean and SD).
547 Results in Table 5 have time interval length $t_{mins} = 1$ minute in order to
548 report the most-detailed feature values. For further change analysis, deci-
549 sion trees are shown in Figure 3 for HS profiles HS1-4. Similarly for the
550 BF dataset, each participant’s change scores and change significance test-
551 ing results are presented in Table 6. The five contextual features (number
552 of bouts, bout minutes, bout steps, daily steps, and sedentary percent) pre
553 and post-intervention values are listed in Table 7. Finally, decision trees
554 are shown in Figure 4 for select BF participants with a significant virtual
555 classifier change score (BF3, BF26, and BF29).

556 5. Discussion

557 We investigate unsupervised change detection and analysis for step-based
558 time series. We compare five change detection approaches, four from the lit-
559 erature and one proposed algorithm. We also implement change significance
560 testing and compute several features for explaining detected changes. The

Table 3: Hybrid-synthetic (HS) significant change detection as a function of time interval size t_{mins} for each HS profile. Results are in the form Count: Boolean (is change significant? 0: false, 1: true) {HS0, HS1, HS2, HS3, HS4}.

t_{mins}	RuSIF	Texture-based	sw-PCAR	Virtual classifier	Total
1	2:0,0,1,0,1	1:0,0,1,0,0	3:0,1,1,0,1	4:0,1,1,1,1	10
5	3:0,1,1,0,1	2:0,0,1,0,1	3:0,1,1,0,1	4:0,1,1,1,1	12
10	2:0,0,1,0,1	2:0,0,1,1,0	2:0,1,1,0,0	3:0,0,1,1,1	9
15	3:0,1,1,0,1	3:0,1,1,1,0	1:0,1,0,0,0	3:0,0,1,1,1	10
20	3:0,1,1,0,1	4:0,1,1,1,1	1:0,1,0,0,0	3:0,0,1,1,1	11
25	1:0,0,1,0,0	2:0,0,1,1,0	1:0,1,0,0,0	3:0,0,1,1,1	7
30	3:0,1,1,0,1	1:0,0,0,1,0	1:0,1,0,0,0	4:0,1,1,1,1	9
35	2:0,1,1,0,0	1:0,0,0,1,0	1:0,1,0,0,0	3:0,1,1,0,1	7
40	3:0,1,1,1,0	0:0,0,0,0,0	1:0,1,0,0,0	3:0,0,1,1,1	7
45	2:0,0,1,0,1	0:0,0,0,0,0	1:0,1,0,0,0	2:0,0,1,0,1	5
50	2:0,0,1,0,1	0:0,0,0,0,0	1:0,1,0,0,0	4:0,1,1,1,1	7
55	3:0,1,1,1,0	1:0,0,0,1,0	0:0,0,0,0,0	3:0,1,1,0,1	7
60	4:0,1,1,1,1	0:0,0,0,0,0	1:0,1,0,0,0	4:0,1,1,1,1	9
Total	33:0,8,13,3,9	17:0,2,6,7,2	17:0,12,3,0,2	43:0,7,13,10,13	110

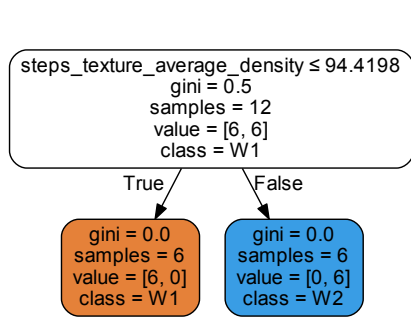
0 = insignificant, 1 = significant

Table 4: Hybrid-synthetic (HS) and B-Fit (BF) significant change detection as a function of time interval size t_{mins} for each HS profile. HS results are in the form PCAR change values {HS0, HS1, HS2, HS3, HS4}. BF results list the top 3 participants with the highest PCAR change scores.

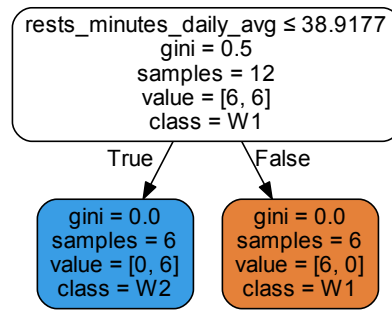
t_{mins}	Hybrid-synthetic PCAR (%)	B-Fit PCAR (%)
1	0.67,48.10,29.89,43.43,2.89	BF27:1.52,BF30:0.73,BF20:0.04
5	32.95,53.33,38.07,44.67,13.77	BF27:3.39,BF30:1.87,BF20:1.62
10	31.22,50.00,52.06,48.00,23.45	BF3:10.04,BF27:7.30,BF26:5.83
15	0.00,50.71,56.34,45.80,10.00	BF3:36.84,BF27:9.02,BF1:4.73
20	0.00,53.92,65.66,61.29,17.31	BF3:34.68,BF27:24.51,BF30:9.47
25	0.00,55.00,66.22,59.46,12.05	BF3:38.95,BF30:11.39,BF23:5.88
30	0.00,70.77,65.15,60.66,27.27	BF3:41.18,BF27:17.91,BF1:7.58
35	0.00,56.36,80.85,46.15,22.95	BF3:52.23,BF27:24.14,BF20:22.22
40	0.00,60.87,58.70,55.32,20.76	BF3:52.83,BF27:9.26,BF29:8.70
45	0.00,65.91,63.42,51.28,25.00	BF3:51.06,BF1:15.23,BF27:8.33
50	22.86,63.16,57.14,50.00,27.5	BF3:48.98,BF27:23.81,BF1:22.73
55	0.00,62.86,63.64,54.55,42.11	BF3:50.00,BF27:31.58,BF1:15.39
60	0.00,70.00,72.41,53.33,33.33	BF3:51.43,BF1:18.75,BF27:14.29

Table 5: Hybrid-synthetic (HS) feature results (mean \pm standard deviation) with $t_{mins} = 1$ minute. Window one and window two values are separated by a comma.

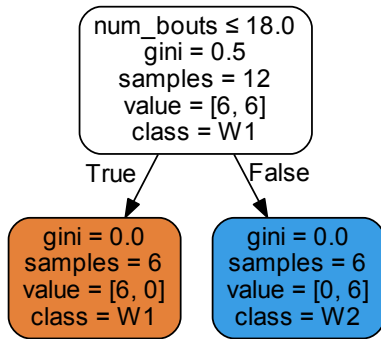
ID	Number of bouts	Bout minutes	Daily steps	Sedentary %
HS0	70.33, 70.00	5.10 \pm 9.91, 5.13 \pm 9.92	20601.65, 21274.32	75.65, 75.56
HS1	34.50, 14.17	19.39 \pm 23.93, 46.82 \pm 56.78	36409.49, 72769.11	64.57, 54.59
HS2	71.50, 62.50	5.07 \pm 9.83, 7.63 \pm 11.13	20755.53, 30037.48	75.62, 67.44
HS3	54.83, 102.83	18.71 \pm 51.74, 6.49 \pm 14.49	14395.85, 14746.43	45.22, 63.72
HS4	53.50, 81.33	8.14 \pm 12.61, 4.56 \pm 8.48	22048.02, 17327.00	70.66, 77.86



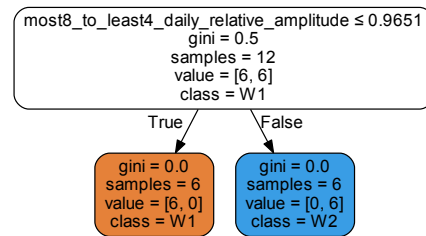
(a) HS1



(b) HS2



(c) HS3



(d) HS4

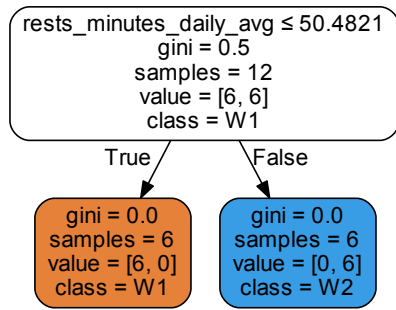
Figure 3: Decision trees for HS profiles with significant virtual classifier change scores for $t_{mins}=5$ minutes.

Table 6: B-Fit (BF) significant change detection as a function of time interval size t_{mins} . Results are in the sparse form Count:IDs: (Boolean is change significant? 0 false, 1 true).

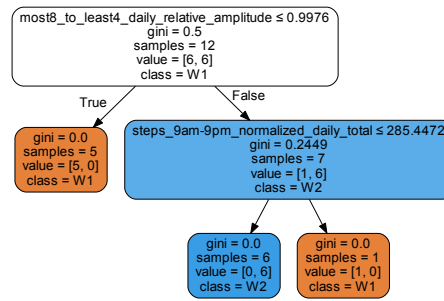
t_{mins}	RulSIF	Texture	sw-PCAR	Virtual classifier	Total
1	1:BF3:1	0	10:BF29:0	5:BF3,24,26,29,30:1	16
5	2:BF3,26:1	0	5:BF1,3,18,20,27:1	6:BF3,24,26,27,29,30:1	13
10	1:BF3:1	0	4:BF1,3,20,27:1	6:BF3,18,24,26,27,29:1	11
15	1:BF3:1	0	4:BF1,3,20,27:1	4:BF3,26,29,30:1	9
20	2:BF3,26:1	0	3:BF3,20,27:1	4:BF3,26,29,30:1	9
25	1:BF3:1	0	3:BF3,20,27:1	2:BF26,29:1	6
30	2:BF3,27:1	0	3:BF3,20,27:1	2:BF24,29:1	7
35	3:BF3,20,27:1	0	3:BF3,20,27:1	5:BF3,24,26,29,30:1	11
40	1:BF3:1	0	3:BF3,20,27:1	3:BF3,26,29:1	7
45	4:BF3,18,20,24	0	1:BF3:1	2:BF3,29:1	7
50	1:BF3:1	0	1:BF3:1	5:BF3,24,26,28,29:1	7
55	1:BF29:1	0	1:BF3:1	3:BF3,20,29:1	5
60	2:BF20,28:1	1:BF27:1	1:BF3:1	4:BF3,24,29,30	8
Total	22	1	42	51	116

Table 7: B-Fit (BF) feature results (mean \pm standard deviation) with $t_{mins} = 1$ minute. Pre and post intervention values are separated by a comma.

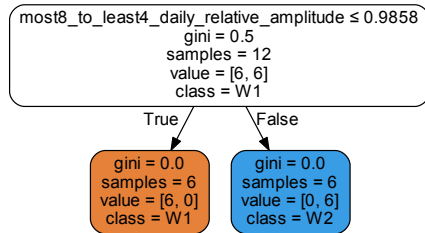
ID	Number of bouts	Bout minutes	Daily steps	Sedentary %
BF1	73.50, 89.67	2.35 \pm 1.93, 2.63 \pm 2.61	3479.00, 4658.33	88.37%, 84.43%
BF3	81.00, 15.83	2.57 \pm 2.54, 2.75 \pm 1.85	4279.50, 1161.44	86.30%, 97.44%
BF18	88.50, 27.50	2.72 \pm 2.86, 2.36 \pm 2.08	5886.67, 4558.50	84.06%, 86.32%
BF20	81.17, 60.33	3.90 \pm 5.27, 3.71 \pm 4.85	11177.00, 7399.67	79.11%, 85.71%
BF23	73.67, 76.50	2.96 \pm 4.06, 2.60 \pm 3.22	6994.17, 5470.50	85.71%, 86.22%
BF24	105.33, 88.00	2.63 \pm 2.46, 2.66 \pm 2.39	7127.00, 6207.67	81.82%, 84.85%
BF26	64.33, 63.50	3.30 \pm 4.20, 3.45 \pm 3.61	7354.67, 6181.17	85.78%, 85.90%
BF27	104.17, 102.50	5.52 \pm 7.51, 3.35 \pm 3.79	17680.78, 11440.00	66.66%, 77.18%
BF28	99.50, 116.67	2.40 \pm 2.49, 2.40 \pm 2.57	5844.50, 6731.83	84.11%, 81.46%
BF29	85.00, 80.50	2.99 \pm 3.24, 3.58 \pm 4.29	1136.51, 1210.85	82.94%, 81.62%
BF30	83.00, 89.50	2.51 \pm 2.73, 2.31 \pm 2.23	5753.50, 4868.83	86.16%, 86.44%



(a) BF3



(b) BF26



(c) BF29

Figure 4: Decision trees for B-Fit (BF) participants with significant virtual classifier change scores for $t_{mins} = 5$ minutes.

561 abilities of the presented methods to detect change are evaluated on two orig-
562 inal datasets: 1) 5 synthetic profiles and 2) 11 participant’s Fitbit data from
563 an intervention study.

564 5.1. Hybrid Synthetic Dataset

565 The HS dataset reveals several insights into the change detection algo-
566 rithms. First, the time interval length yielding the highest number of signif-
567 icant changes is $t_{mins} = 5$ minutes with 12 changes, closely followed by t_{mins}
568 $= 20$ minutes with 11 changes (see Table 3). Since HS profiles are sampled
569 from a volunteers real user Fitbit data, these intervals suggest movement
570 patterns occur in 5 and 20 minute chunks for this individual. For all time
571 interval lengths, the algorithms do not detect a significant change between
572 window one and window two data for the HS0 profile. HS0 is generated to
573 exhibit small day-to-day variation in step intensity and is not characterized
574 by large changes between windows; however, PCAR detects an average of
575 6.75% change for HS0. We can use PCAR HS0 values as baseline change
576 scores for relative interpretations of PCAR HS1-4 values.

577 HS2 and HS4 were generated to exhibit abrupt changes between the first
578 and second window, whereas HS1 and HS3 exhibit gradual day-to-day change
579 from day one of window one to the last day of window 2. For HS1-4 pro-
580 files, significant changes between windows are detected. For all time inter-
581 val lengths, the virtual classifier approach picks up the most changes (43
582 changes), followed by RuLSIF (33), sw-PCAR (17), and texture-based (17).
583 Changes for HS2 (35) and HS1 (29) are the most frequently detected, fol-
584 lowed by HS4 (26) and HS3 (20). As a group, the algorithms’ are able to sense
585 change in value (HS1, HS2) and changes in variability (HS3, HS4), with 64
586 and 46 changes respectively. RuLSIF struggles to detect gradual variability
587 change (HS3: 8), but perfectly detects window-based value change (HS2: 13)
588 for all time intervals. In fact, perfect detections are made by virtual classifier
589 (HS2, HS4: 13) and near perfect detections are made for sw-PCAR (HS1:
590 12). Investigating the $t_{mins} = 5$ minutes results reveal all four algorithms
591 determine significant changes for HS2 and HS4 (see Table 3); however, for
592 HS2 and HS4, PCAR identifies lower change (38.07% and 13.77%) than the
593 gradual change profiles HS1 and HS3 (53.33% and 44.67%). This HS4 PCAR
594 score is less than the $t_{5minutes}$ HS0 PCAR score, 32.95%, implying PCAR did
595 not detect noteworthy change for HS4.

596 Upon inspection of the associated decision trees for HS2-4 (see Figure 3,
597 the features of texture density, average daily rest minutes, number of bouts,

598 and window-based relative amplitude are discriminatory features. The ex-
599 planatory power of the features is potentially useful for reporting to the wear-
600 able sensor user the dimensions of change in their physical activity. Features
601 useful for such purposes are simple, common features that do not require
602 interpretation. For example, texture density or relative amplitude are useful
603 features for detecting changes in PA patterns, but are relatively unimportant
604 to a user. More meaningful features to a user include number of bouts, num-
605 ber of steps per bout, rest period minutes, and sedentary percent. Table 5
606 shows these features for the HS profiles. HS0 exhibits quite similar window
607 one and window two values for all features. HS2 and HS4 both have small
608 standard deviations due window-based change in lieu of day-to-day change
609 (HS1 and HS3).

610 5.2. *B-Fit Dataset*

611 Analyzing the BF participants' data poses additional challenges that are
612 not present with the HS profiles. Real-world human subject data is inher-
613 ently noisy, characterized by seemingly random bouts of PA and rest periods.
614 Furthermore, self-report and direct measurement of physical activity are of-
615 ten not congruent, with previous studies reporting correlations in as wide
616 of a range of -0.71 to 0.96 [26]. For the BF group, Table 2 shows a wide
617 spread of self-reported goal achievement ratings for the exercise and cardio-
618 vascular risk factor categories. For example, BF24, BF28, and BF29 rated
619 their exercise goal achievements low (exercise: 0.6, 0, 1 respectively). Due
620 to heart problems, BF28's doctor instructed him not participate in exercise-
621 related activities. On the other hand, BF3 rated their goal achievements
622 the highest (exercise: 3; cardio 2.86). Upon inspection of BF3's data, it is
623 evident there is a discrepancy between the participant's perception of her
624 PA and the steps recorded by the Fitbit (specifically 9/14-9/19 compared to
625 11/19-11/21). It is not uncommon for self-reported measures of physical ac-
626 tivity to be inconsistent with direct measures [26]; therefore, the self-ratings
627 presented in Table 2 are used for insights into individual goal achievements,
628 not as ground truth information for changes exhibited. The issues with self-
629 reported PA measures exacerbate the need for unsupervised change detection
630 and analysis methods.

631 Depending on the algorithm, significant changes are commonly detected
632 for 5 out of the 11 BF participants: BF3: 35; BF27: 14; BF29:14; BF20: 13,
633 BF26: 12; (see Table 3). Virtual classifier and sw-PCAR detect the highest
634 number of changes (51 and 42 changes each), but the distribution of detected

635 changes is highly influenced by time interval length (3.92 ± 1.44 , 3.23 ± 2.42
636 number of changes detected respectively). sw-PCAR is not sensitive for small
637 time intervals ($t_{mins}=1$ minutes) or large time intervals ($t_{mins}=\{45, 50, 55,$
638 60 minutes}), and the number of changes detected decreases as time interval
639 length increases. Virtual classifier does not appear to be as heavily influenced
640 by the time interval length. The texture-based approach is the least sensitive
641 algorithm, only detecting change for BF27 with $t_{mins} = 60$ minutes. Finally,
642 PCAR detects the most changes with time window $t_{mins} = 35$ minutes.

643 Performing change analysis and investigating the changes detected yields
644 insights for several of the participants. BF3 rated herself as completely meet-
645 ing her exercise goal of walking more; however, the Fitbit data tells a different
646 story. Several features in Table 5 show decreased PA for BF3: average num-
647 ber of bouts (W1: 81.00, W2: 15.83), daily steps (W1: 4279.50, W2: 1161.44
648 steps), and percentage of time sedentary (W1: 86.30%, W2: 97.44%). Addi-
649 tionally, BF3’s decision tree (see Figure 4a) provides evidence that she rested
650 more during post-intervention testing. In summary, the features suggest the
651 changes detected by the algorithms are actually changes in the opposite di-
652 rection of her goal. Contrary to BF3, BF29 exhibited a significant change
653 (as detected consistently by virtual classifier) in the direction towards her
654 goal of walking more. Inspection of BF29’s features shows an increase bout
655 minutes, bout steps, and average steps per day. Average daily steps increased
656 from 1136.51 steps pre-intervention to 1210.85 steps post-intervention test-
657 ing, a 6.54% increase. The remaining participants with significant changes
658 (BF3, BF20, BF24, BF26, BF27, and BF30) demonstrated a decrease in av-
659 erage daily steps taken from pre to post intervention. While this suggests
660 the exercise component of the intervention was not successful for these par-
661 ticipants, the participants’ physical activity levels may have been influenced
662 by seasonal effects [27]. Pre-intervention Fitbit data collection occurred
663 in September, which is considerably warmer than November/December in
664 Washington State, which was the period of post-intervention data collection.
665 It is also worth noting the participants exhibited improvements in other phys-
666 ical activity features. For example, relative amplitude has been reported to
667 decrease with worsening health [28], thus BF26 and BF29’s increased rela-
668 tive amplitude post-intervention is healthy (see Figures 4b and 4c). Also,
669 BF28 was not planning on increasing exercise; however, BF28 increased their
670 daily steps post-intervention by 15.18%.

671 One of the limitations of this study includes only having one week of
672 pre-intervention Fitbit data for BF participants. With at least two weeks

673 of pre-intervention data, change scores can be computed between week one
674 and two of pre-intervention data to provide an estimate of inter-week vari-
675 ability. With a quantification of inter-week variability, we can determine
676 if the change measured between pre and post-intervention weeks is due to
677 the intervention or natural variability. An additional limitation includes not
678 having full 7 days of BF data during pre and post intervention weeks. Fi-
679 nally, more sophisticated methods to fill missing data could be utilized with
680 fitness trackers that include heart rate monitors, due to more reliable detec-
681 tion of sensor donned/doffed. Consequently, future work includes performing
682 change analysis on real-world datasets from different fitness trackers, mul-
683 tivariate data (i.e. heart rate, elevation, etc.), labeled activity data, and
684 longer windows of time. With time series data longer than two years, sev-
685 eral additional analyses could be performed: daily/weekly/monthly/yearly
686 period analysis, slicing along different dimensions (i.e. Mondays, weekends,
687 holidays, or activities if labeled information is available), etc.

688 **6. Conclusion**

689 We address the problem of unsupervised physical activity change detec-
690 tion and analysis. We compare five change detection approaches, four from
691 the literature and one proposed algorithm. We objectively compare the algo-
692 rithms' abilities to capture different types of changes in five distinct synthetic
693 datasets representing realistic changes in physical activity patterns. We also
694 evaluate the algorithms on real-world physical activity data collected from an
695 intervention study where 11 Fitbit users set goals to improve various facets of
696 their health. The results indicate the algorithms detected the most significant
697 changes in both datasets for time interval lengths of 1, 5, and 15 minutes. For
698 the synthetic dataset, the virtual classifier and sw-PCAR approaches picked
699 up on the highest number of changes. Changes for profiles exhibiting large
700 changes between windows are more likely to be detected than those exhibiting
701 incremental day-to-day changes. For the real world dataset, the algorithms
702 frequently detect changes for 5 of the 11 participants. Change analysis for
703 these 5 participants' physical activity features reveal only 1 exhibits an in-
704 crease in daily steps taken post-intervention. Contextual features such as
705 average number of daily steps, minutes spent resting, number of steps per
706 bout, and sedentary percent provide an explanation of the changes detected.
707 The algorithms and analysis methods are useful data mining techniques for
708 unsupervised, window-based change detection. Future work involves imple-

709 menting the algorithms in an online, smartphone application to track users'
710 physical activity and motivate their progress towards their health goals.

711 **Acknowledgements**

712 We wish to thank Catherine Sumida and Thao Vo for their help in data
713 collection. This work is supported in part by [Grant?].

714 **References**

- 715 [1] P. Lally, C. van Jaarsveld, H. Potts, J. Wardle, How are habits formed:
716 Modelling habit formation in the real world, *Eur. J. Soc. Psychol.* 40 (6)
717 (2010) 998–1009.
- 718 [2] B. H. Dobkin, Wearable motion sensors to continuously measure real-
719 world physical activities, *Curr Opin Neurol* 26 (6) (2013) 602–608.
- 720 [3] J. Merilahti, P. Viramo, I. Korhonen, Wearable monitoring of physical
721 functioning and disability changes, circadian rhythms and sleep patterns
722 in nursing home residents, *IEEE Journal of Biomedical and Health In-*
723 *formatics PP* (99) (2015) 1–1.
- 724 [4] P. Paavilainen, I. Korhonen, J. Ltjnen, L. Cluitmans, M. Jylh, A. Srel,
725 M. Partinen, Circadian activity rhythm in demented and non-demented
726 nursing-home residents measured by telemetric actigraphy, *J Sleep Res*
727 14 (1) (2005) 61–68.
- 728 [5] S. Wang, M. Skubic, Y. Zhu, Activity Density Map Visualization and
729 Dissimilarity Comparison for Eldercare Monitoring, *IEEE Transactions*
730 *on Information Technology in Biomedicine* 16 (4) (2012) 607–614.
- 731 [6] T.-H. Tan, M. Gochoo, K.-H. Chen, F.-R. Jean, Y.-F. Chen, F.-J. Shih,
732 C. F. Ho, Indoor activity monitoring system for elderly using RFID and
733 Fitbit Flex wristband, in: *2014 IEEE-EMBS International Conference*
734 *on Biomedical and Health Informatics (BHI)*, 2014, pp. 41–44.
- 735 [7] P. N. Dawadi, D. J. Cook, M. Schmitter-Edgecombe, Modeling patterns
736 of activities using activity curves, *Pervasive and Mobile Computing*

- 737 [8] S. Liu, M. Yamada, N. Collier, M. Sugiyama, Change-point detection in
738 time-series data by relative density-ratio estimation, *Neural Networks*
739 43 (2013) 72–83.
- 740 [9] S. Hido, T. Id, H. Kashima, H. Kubo, H. Matsuzawa, Unsupervised
741 Change Analysis Using Supervised Learning, in: T. Washio, E. Suzuki,
742 K. M. Ting, A. Inokuchi (Eds.), *Advances in Knowledge Discovery and*
743 *Data Mining*, no. 5012 in *Lecture Notes in Computer Science*, Springer
744 Berlin Heidelberg, 2008, pp. 148–159, doi: 10.1007/978-3-540-68125-
745 0_15.
- 746 [10] M. Yamada, A. Kimura, F. Naya, H. Sawada, Change-point Detection
747 with Feature Selection in High-dimensional Time-series Data, in: *Pro-*
748 *ceedings of the Twenty-Third International Joint Conference on Arti-*
749 *ficial Intelligence, IJCAI '13*, AAAI Press, Beijing, China, 2013, pp.
750 1827–1833.
- 751 [11] D.-H. Tran, M. M. Gaber, K.-U. Sattler, Change Detection in Streaming
752 Data in the Era of Big Data: Models and Issues, *SIGKDD Explor.*
753 *Newsl.* 16 (1) (2014) 30–38.
- 754 [12] K. Feuz, D. Cook, C. Rosasco, K. Robertson, M. Schmitter-Edgecombe,
755 Automated Detection of Activity Transitions for Prompting, *IEEE*
756 *Transactions on Human-Machine Systems* 45 (5) (2015) 575–585.
- 757 [13] F. Javed, S. Farrugia, M. Colefax, K. Schindhelm, Early Warning of
758 Acute Decompensation in Heart Failure Patients Using a Noncontact
759 Measure of Stability Index, *IEEE Transactions on Biomedical Engineer-*
760 *ing* 63 (2) (2016) 438–448.
- 761 [14] M. Hu, S. Zhou, J. Wei, Y. Deng, W. Qu, Change-point Detection in
762 Multivariate Time-series Data by Recurrence Plot, *WSEAS Transac-*
763 *tions on Computers* 13 (2014) 592–599.
- 764 [15] Y. Xu, Z. Zhang, P. Yu, B. Long, Pattern Change Discovery Between
765 High Dimensional Data Sets, in: *Proceedings of the 20th ACM Interna-*
766 *tional Conference on Information and Knowledge Management, CIKM*
767 *'11*, ACM, New York, NY, USA, 2011, pp. 1097–1106.

- 768 [16] W. Ng, M. Dash, A change detector for mining frequent patterns over
769 evolving data streams, in: IEEE International Conference on Systems,
770 Man and Cybernetics, 2008. SMC 2008, 2008, pp. 2407–2412.
- 771 [17] C. J. Caspersen, K. E. Powell, G. M. Christenson, Physical activity,
772 exercise, and physical fitness: definitions and distinctions for health-
773 related research., *Public Health Rep* 100 (2) (1985) 126–131.
- 774 [18] F. Albrechtsen, others, Statistical texture measures computed from gray
775 level cooccurrence matrices, Image processing laboratory, Department of
776 informatics, university of oslo (2008) 1–14.
- 777 [19] L. Chen, J. Hoey, C. Nugent, D. Cook, Z. Yu, Sensor-Based Activity
778 Recognition, *IEEE Transactions on Systems, Man, and Cybernetics,*
779 *Part C: Applications and Reviews* 42 (6) (2012) 790–808.
- 780 [20] Y. Benjamini, Y. Hochberg, Controlling the False Discovery Rate: A
781 Practical and Powerful Approach to Multiple Testing, *Journal of the*
782 *Royal Statistical Society. Series B (Methodological)* 57 (1) (1995) 289–
783 300.
- 784 [21] O. Maimon, L. Rokach, *Data mining and knowledge discovery hand-*
785 *book*, Springer, New York, 2005.
- 786 [22] J. W. Tukey, *Exploratory data analysis*.
- 787 [23] F. E. Grubbs, Procedures for Detecting Outlying Observations in Sam-
788 ples, *Technometrics* 11 (1) (1969) 1.
- 789 [24] V. J. Hodge, J. Austin, A Survey of Outlier Detection Methodologies,
790 *Artif Intell Rev* 22 (2) (2004) 85–126.
- 791 [25] R. Refinetti, G. C. Lissen, F. Halberg, Procedures for numerical analysis
792 of circadian rhythms, *Biol Rhythm Res* 38 (4) (2007) 275–325.
- 793 [26] S. A. Prince, K. B. Adamo, M. E. Hamel, J. Hardt, S. C. Gorber,
794 M. Tremblay, A comparison of direct versus self-report measures for
795 assessing physical activity in adults: a systematic review, *International*
796 *Journal of Behavioral Nutrition and Physical Activity* 5 (2008) 56.
- 797 [27] P. Tucker, J. Gilliland, The effect of season and weather on physical
798 activity: A systematic review, *Public Health* 121 (12) (2007) 909–922.

799 [28] J. Merilahti, T. Petkoski-Hult, M. Ermes, M. v. Gils, H. Lahti, A. Yli-
800 nen, L. Autio, E. Hyvrinen, J. Hyttinen, Evaluation of new concept for
801 balance and gait analysis: patients with neurological disease, elderly
802 people and young people, *Gerontechnology* 7 (2) (2008) 164.