# Automatic physical activity and in-vehicle status classification based on GPS and accelerometer data: A hierarchical classification approach using machine learning techniques

**2 authors:**

Kangjae Lee
University of Illinois, Urbana-Champaign
**10** PUBLICATIONS **38** CITATIONS

SEE PROFILE

Mei-Po Kwan
University of Illinois, Urbana-Champaign
**230** PUBLICATIONS **8,757** CITATIONS

SEE PROFILE

**Some of the authors of this publication are also working on these related projects:**

311 Non-Emergency Municipal Services View project

Physical activity classification using accelerometer and GPS data View project

RESEARCH ARTICLE

WILEY Transactions in GIS

# Automatic physical activity and in-vehicle status classification based on GPS and accelerometer data: A hierarchical classification approach using machine learning techniques

Kangjae Lee [iD] | Mei-Po Kwan [iD]

University of Illinois at Urbana-Champaign, Urbana, Illinois

**Correspondence**
Kangjae Lee, Illinois Informatics Institute, University of Illinois at Urbana-Champaign, 1205 West Clark, Urbana, IL 61801, USA.
Email: kasbiss@gmail.com

## Abstract

Due to the advancement of tracking technology, a large quantity of movement data has been collected and analyzed in various research domains. In human mobility and physical activity (PA) research, GPS trajectories and the capabilities of geographic information systems (GIS) facilitate a better understanding of the associations between PA and various environmental factors taking individuals' daily travels into account. PA research, however, needs to widen its focus from the intensity of PA to types of PA, which may provide useful clues for understanding specific health behaviors in particular geographic contexts. This study proposes and develops an algorithm to automatically classify PA types and in-vehicle status using GPS and accelerometer data. Walking, standing, jogging, biking and sedentary/in-vehicle statuses are identified through hierarchical classification processes based on machine learning and geospatial techniques. The proposed algorithm achieved high predictive accuracy on real-world GPS and accelerometer data. It can greatly reduce participants' and researchers' burdens by automatically identifying PA types and in-vehicle status for human mobility research, which is also known as travel mode imputation in transportation research.

# 1 | INTRODUCTION

Due to the advancement of tracking technology, a large quantity of movement data has been collected and analyzed in various research domains (Eagle, Pentland, & Lazer, 2009; Gonzalez, Hidalgo, & Barabasi, 2008; Shoval & Isaacson, 2007). As one of the tracking devices, global positioning system (GPS) receivers have been widely used to collect data that enhance our understanding of the spatiotemporal dynamics of moving objects, such as animals (Dodge et al., 2013; Laube, Dennis, Forer, & Walker, 2007), humans (Kwan, 2004; Shoval et al., 2011; Wang, Kwan, & Chai, 2018), or vehicles (Downs & Horner, 2012; Ferreira, Poco, Vo, Freire, & Silva, 2013). In physical activity (PA) research, GPS trajectories and the capabilities of geographic information systems (GIS) facilitate a better understanding of the associations between moderate to vigorous physical activity (MVPA), such as brisk walking and running, and various environmental factors, taking individuals' daily travels into account (Almanza, Jerrett, Dunton, Seto, & Pentz, 2012; Boruff, Nathan, & Nijnstein, 2012; Cooper et al., 2010; Jansen, Ettema, Pierik, & Dijst, 2016; Lachowycz, Jones, Page, Wheeler, & Cooper, 2012; Rodríguez et al., 2012; Troped, Wilson, Matthews, Cromley, & Melly, 2010). Such recent PA studies that use GPS trajectories reveal the importance of non-residential areas in people's health behaviors and outcomes, in addition to residential neighborhoods (Diez Roux & Mair, 2010; Perchoux, Chaix, Cummins, & Kestens, 2013). High-resolution GPS data can also be used to mitigate the uncertain geographic context problem (UGCoP), since they help identify the various non-residential contexts that may affect people's health (Kwan, 2012a, b, 2013).

PA research has largely used objectively measured accelerometer data to assess the intensity of people's PA, like MVPA, since objective PA measures yield more significant findings than subjective and self-reported measures (Browning & Lee, 2017). PA research, however, needs to widen its focus from the intensity of PA to types of PA, which may provide useful clues for understanding specific health behaviors in particular geographic contexts (Jankowska, Schipperijn, & Kerr, 2015). Some PA conceptual models suggest that specific types of PA, like walking and biking, have associations with specific environmental or geographic contexts, such as trails, safe pedestrian sidewalks, and supportive facilities. Yet, more consistent empirical evidence is needed to support such associations (Loukaitou-Sideris, 2006; Sallis, Bauman, & Pratt, 1998). In addition, few PA studies to date using GPS and accelerometer data have taken into account motorized-transport modes like traveling in a vehicle (e.g., Voss et al., 2016). According to Gordon-Larsen, Nelson, and Beam (2005), the high utilization of non-motorized transport (e.g., walking and biking) for commuting is mostly exhibited among the group of young adults who meet PA recommendations, whereas the higher percentage of young adults who do not meet PA recommendations is associated with higher utilization of motorized-transport modes (e.g., public transit and private car) for traveling to workplaces or schools. Further, the impact of various environmental factors on people's PA may also be less when they stay inside a vehicle, compared to when they walk or run outside. It is thus important to identify and separate motorized-transport modes (vehicle-based movement) from PA (non-vehicle-based human movement) in order to more accurately estimate people's exposure to various environmental contexts. However, only a few PA studies have attempted to identify whether people are traveling in vehicles (in-vehicle status) or not, through the automatic classification of travel modes and PA to date (Ellis et al., 2014; Zhou, 2014).

To contribute to this literature, this study proposes and develops an algorithm to automatically classify PA types and in-vehicle status using GPS and accelerometer data available to the public. Hierarchical classification processes based on machine learning techniques are innovative approaches adopted in this study. As a branch of artificial intelligence, machine learning improves the prediction of outcomes through a training process based on machine learning models/algorithms using a large amount of input data. The introduction of hierarchical classification is to jointly identify more classes (identified with distinctive labels)—like the different types of PA in this study—using heterogeneous sensor data, such as GPS and accelerometer data, than using only one of these. The hierarchical classification processes are discussed in detail in Section 3. In this study, three components constitute the framework of the hierarchical classification algorithm: indoor/outdoor classification, classification using GPS data (outdoors), and classification using accelerometer data (indoors). Machine learning techniques make

predictions of people's PA types and in-vehicle status based on the generated numeric or categorical features from collected GPS and accelerometer data. As measurable characteristics, features are informative quantifiable properties calculated from input data and one of the fundamental elements in machine learning used to predict different classes (e.g., PA types in this study). Regarding predicted classes, four PA types (biking, running, walking, standing), sitting (sedentary status), and traveling in a vehicle (in-vehicle status) are automatically identified through the developed algorithm. Because running and biking are the two most popular outdoor PA types for young adults in the US (Outdoor Foundation, 2016), the proposed algorithm will seek to identify these two PA types.

The proposed algorithm is quantitatively and qualitatively validated using real-world GPS data collected from three subjects in highly and moderately urbanized areas. Highly urbanized areas here mean areas with high building density, tall buildings, and heavy traffic, while moderately urbanized areas are areas with moderate building density, fewer tall buildings, and moderate traffic. The rationale for choosing highly and moderately urbanized areas in this study is that the degree of urbanization (e.g., tall buildings and traffic congestion) may influence the accuracy of GPS measurement. For example, in highly urbanized areas, GPS positional errors are likely to be higher than those in moderately urbanized areas due to the obstruction of GPS signals by tall buildings (known as urban canyons in GPS parlance). Further, travel speed tends to be very slow in certain highly congested road segments in highly urbanized areas, which may affect the classification accuracy of in-vehicle status, since the GPS position of a moving object becomes unstable when it is stationary or moving slowly. Therefore, the performance of the algorithm needs to be explored in these two types of areas with different levels of urbanization.

In this study, Chicago was chosen as an example of a highly urbanized area, and Champaign, IL was selected as an example of a moderately urbanized area. Chicago is the second largest city in the US, with the most skyscrapers 490 ft (150 m) or greater in height, whereas Champaign has no skyscrapers (Council on Tall Buildings and Urban Habitat, 2018). Chicago also ranked third in total travel delay and total congestion cost among the largest urban areas in the US in 2014 (Schrank, Eisele, Lomax, & Bak, 2015). Regarding total travel delay, Champaign had 1,966,000 hours of extra travel time in 2014, which is much less than Chicago (302,609,000 hours).

The algorithm greatly reduces participants' burdens in recording travel modes by automatically classifying PA types and in-vehicle status, and is thus useful for human mobility research using GPS, including transportation and public health studies. Researchers can also benefit from the automatic classification, since there is no need to be concerned about missing or inaccurate travel mode records made by participants. Automatic travel mode imputation will also be useful for identifying travel patterns associated with particular route choices and for understanding the underlying decision-making processes in route choice research (Broach, Dill, & Gliebe, 2012; Papinski, Scott, & Doherty, 2009). In addition, this study shows how published heterogeneous sensor data (GPS and accelerometer data) can work together to provide accurate and automatic PA classification using machine learning and geospatial techniques.

The remainder of this article is organized into four sections. Section 2 describes previous studies on PA types and transport-mode classification using both GPS and accelerometer data. Section 3 describes the algorithm with respect to its components, focusing more on indoor/outdoor classification and classification using GPS data. Section 4 validates the proposed algorithm, and Section 5 discusses the research findings and provides conclusions.

## 2 | PAST STUDIES ON PA AND TRANSPORT-MODE CLASSIFICATION

Many PA recognition studies have been conducted, leveraging the potential of the sensing capabilities of smartphones with a particular focus on health promotion and management. Particularly, the accelerometer sensors commonly equipped in most smartphones brought profound enhancement to the automatic recognition of PA, including running, walking, and sitting, with high precision using machine learning techniques (Anguita, Ghio, Oneto, Parra, & Reyes-Ortiz, 2012; Arif, Bilal, Kattan, & Ahamed, 2014; Kwapisz, Weiss, & Moore, 2011; Weiss et

al., 2016; Zhang, McCullagh, Nugent, & Zheng, 2010). For instance, PA classification was employed in a web-based application to monitor the PA of children, obese people, or the elderly, and to encourage them to perform sufficient PA in their daily lives (Weiss et al., 2016).

The automatic recognition of different travel modes using GPS trajectories has drawn the attention of transport researchers, due to the low response rates and high incompletion rates of paper- or phone-based travel surveys. For instance, Zheng, Chen, Li, Xie, and Ma (2010) proposed an approach to classify four kinds of travel modes—walking, driving, traveling by bus, and biking—by segmenting GPS trajectories and extracting features like maximum velocity and acceleration in each GPS segment to understand human mobility and display it on web-based mapping applications (Table 1). The study suggests a way of dividing a GPS trajectory into walking and non-walking segments, stressing the importance of the walking mode as a transition to a non-walking mode (e.g., car, bus, or train). With the segmentation approach and the use of a decision tree, the study achieved a predictive accuracy of over 75%. Along with GPS data, GIS data were employed to achieve higher accuracy in travel mode classification. In most cases, transport network data and data on the related infrastructure (e.g., stops, stations, entrances) have been used (Biljecki, Ledoux, & Van Oosterom, 2013; Gong, Chen, Bialostozky, & Lawson, 2012; Witayangkurn, Horanont, Ono, Sekimoto, & Shibasaki, 2013). For example, Biljecki et al. (2013) extracted road, railway, bus, and tram networks, the locations of bus stops, and train stations from OpenStreetMap to detect 10 transport modes using a fuzzy logic approach, and the algorithm achieved 92% accuracy. Besides transport links, Moiseeva, Jessurun, and Timmermans (2010) proposed a travel mode inference system for the classification of seven transport modes using land-use data. The land-use data helped to increase the predictive accuracy for particular kinds of transport modes that may occur on certain types of land use (e.g., railroad tracks). The transport-mode classification using a Bayesian belief network model achieved 95% accuracy. On the other hand, some studies did not use GIS data because these might not be available for many study areas. However, some of the studies were able to achieve moderate to high predictive accuracy (Xiao, Wang, Fu, & Wu, 2017; Zhu et al., 2016).

With respect to performance improvement, how features are extracted can significantly affect the results of travel-mode classification. For the last decade, some studies considered the focal characteristics of sub-segments for each trip, captured through moving windows sliding on GPS trajectories (Bolbol, Cheng, Tsapakis, & Haworth, 2012; Dodge, Weibel, & Forootan, 2009; van Dijk, 2018; Xiao et al., 2017). Bolbol et al. (2012) applied a fixed-size moving window sliding on speed and acceleration values of multi-segment GPS instances, and the classification of six travel modes using support vector machine (SVM) achieved an accuracy of 88%. Dodge et al. (2009) and Xiao et al. (2017) especially adopted the focal characteristic of GPS trajectories by calculating movement parameters from the GPS points that fall within a sliding window, achieving a predictive accuracy of 82% and 91%, respectively. van Dijk (2018) achieved over 99% predictive accuracy in classifying trips (moving) and activities (staying) by introducing moving spatial and temporal windows.

The combined use of multiple sensor data has recently been introduced to PA and travel-mode classification research due to the increasing availability of various sensors. Patrick et al. (2008) developed the Physical Activity and Location Measurement System (PALMS) to understand PA-related energy expenditure associated with time and space in exposure biology research by incorporating GPS, accelerometer, and heart-rate monitoring sensors. With regard to travel mode identification, accelerometer, magnetometer, and gyroscope data recorded using smartphones enhanced the performance of travel-mode detection (Ellis et al., 2014; Feng & Timmermans, 2013; Fang et al., 2016; Shafique & Hato, 2016; Zhou, 2014). Ellis et al. (2014) particularly compared the predictive accuracy of several machine learning models for the classification of six travel modes, including bus, car, sitting, and walking, based on GPS and accelerometer data collected from two trained assistants in different built-environment settings. Evenson and Furberg (2017) developed a smartphone application for users and researchers to automatically predict PA types using GPS and accelerometer data.

In this study, an automatic PA and in-vehicle status classification algorithm is developed using two kinds of open, publicly available sensor data—GPS and accelerometer data—collected by other researchers (these data will be described in Section 3). Being sedentary is an important stationary behavior, which many PA studies also

**TABLE 1** Characteristic properties of past studies on PA types and transport modes

| Author | Predictive accuracy | Classification algorithm or system | Sensor | Classified activities | Observation unit | GIS data use | Moving window |
|---|---|---|---|---|---|---|---|
| Zheng et al. (2010) | 75.60% | Tree-based model | GPS | Walking, biking, car, bus | Segment | No | No |
| Biljecki et al. (2013) | 91.60% | Fuzzy expert system | GPS | Walking, biking, car, bus, train, tram, underground, ferry, sailing boat, aircraft | Segment | Yes | No |
| Moiseeva et al. (2010) | 95.40% | Bayesian belief network | GPS | Walking, running, biking, motorbike, car, bus, train | Segment | Yes | No |
| Xiao et al. (2017) | 90.77% | XGBoost | GPS | Walking, biking, car, bus and taxi, subway, train | Segment | No | Yes |
| Zhu et al. (2016) | 91.44% | Random forest | GPS | Walking, biking, car, bus | Segment | Yes | No |
| Bolbol et al. (2012) | 88.00% | SVM | GPS | Walking, biking, car, bus, train, underground | Segment | No | Yes |
| Dodge et al. (2009) | 82.00% | SVM | GPS | Pedestrian, biking, car, motorbike | Segment | No | Yes |
| van Dijk (2018) | 99.40% | SVM and random forest | GPS | Move, stay | Point | Yes | Yes |
| Ellis et al. (2014) | 91.90% | Random forest | GPS, accelerometer | Walking, biking, car, bus, sitting, standing | Point | No | Yes |
| Fang et al. (2016) | 86.94% | SVM | Accelerometer, magnetometer, gyroscope | Walking, running, biking, vehicle, staying | Segment | No | Yes |
| Feng and Timmermans (2013) | 85% | Bayesian belief network | GPS, accelerometer | Walking, running, biking, motorbike, car, bus, tram, subway | Point | No | Yes |
| Shafique and Hato (2016) | 99.96% | Random forest | Accelerometer | Walking, biking, car, bus, train, subway | Segment | No | Yes |
| Zhou (2014) | Over 80% | Random forest | GPS, accelerometer | Walking, running, biking, in-vehicle, staying | Point | No | No |

examined in addition to MVPA. The contribution of this study lies in how to utilize two different kinds of sensor data to identify travel modes. As adopted in the research by Dodge et al. (2009) and Xiao et al. (2017), the merit of the GPS trajectory operators at the focal level used in this study is that they can help capture distinctive variations in movement derivatives (e.g., velocity, acceleration) over time for different travel modes. This study, in particular, uses GPS trajectory operators extensively at the focal level with different sizes of moving spatial and temporal windows, as described in Section 3.2. Compared to the research by van Dijk (2018) that used spatial and temporal moving windows (as shown in Table 1), this study extracts more features to maximize the benefit of spatial and temporal moving windows for classifying more travel modes by utilizing two different kinds of open sensor data.

# 3 | A FRAMEWORK OF AUTOMATIC CLASSIFICATION OF PA AND IN-VEHICLE STATUS USING PUBLICLY AVAILABLE GPS AND ACCELEROMETER DATA

The automatic PA and in-vehicle status classification in this study is implemented through hierarchical classification processes to jointly identify a total of six classes: biking, running, walking, standing, being sedentary, and riding in a vehicle (Figures 1 and 2). Hierarchical classification has been used to deal with hierarchical structures in real-world systems regarding classification from the top level to lower levels (Dumais & Chen, 2000; McNamara, Crossley, Roscoe, Allen, & Dai, 2015). In this study, a hierarchical classification approach is expected to greatly facilitate the accurate classification of indoor versus outdoor activities at the top level, and vehicle-based versus non-vehicle-based travel modes using GPS and accelerometer data at the lower levels (Figure 1). That is because GPS data are suitable for identifying vehicle-based movement versus non-vehicle-based human movement,
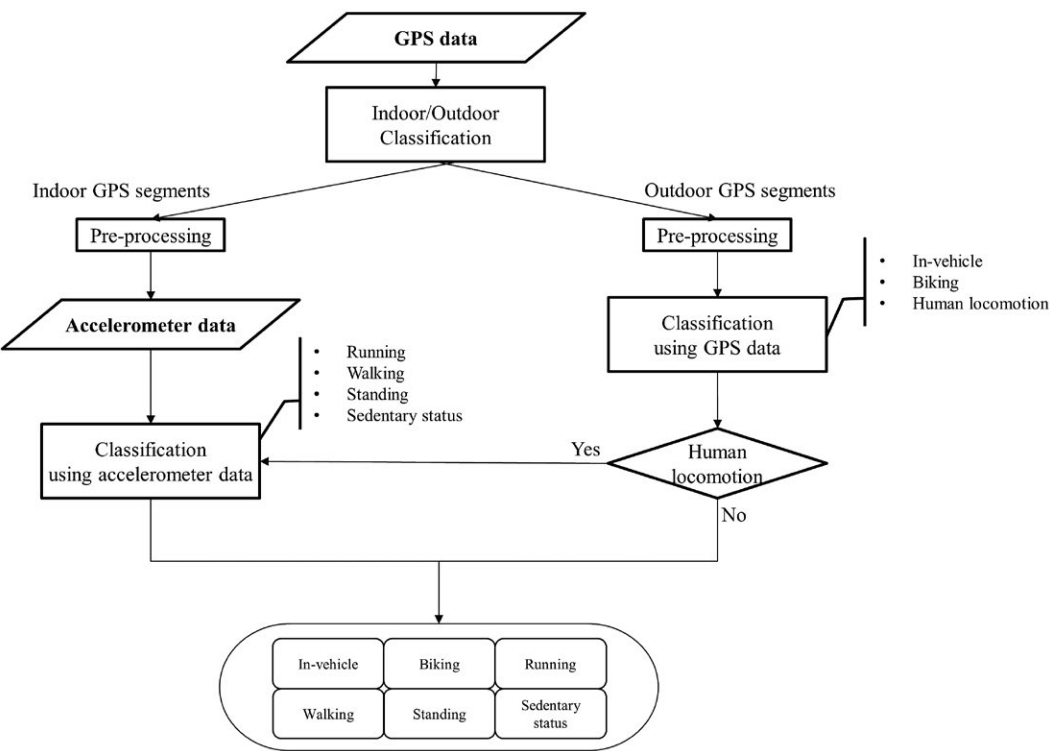


**FIGURE 1** Flow diagram of the PA and in-vehicle classification algorithm
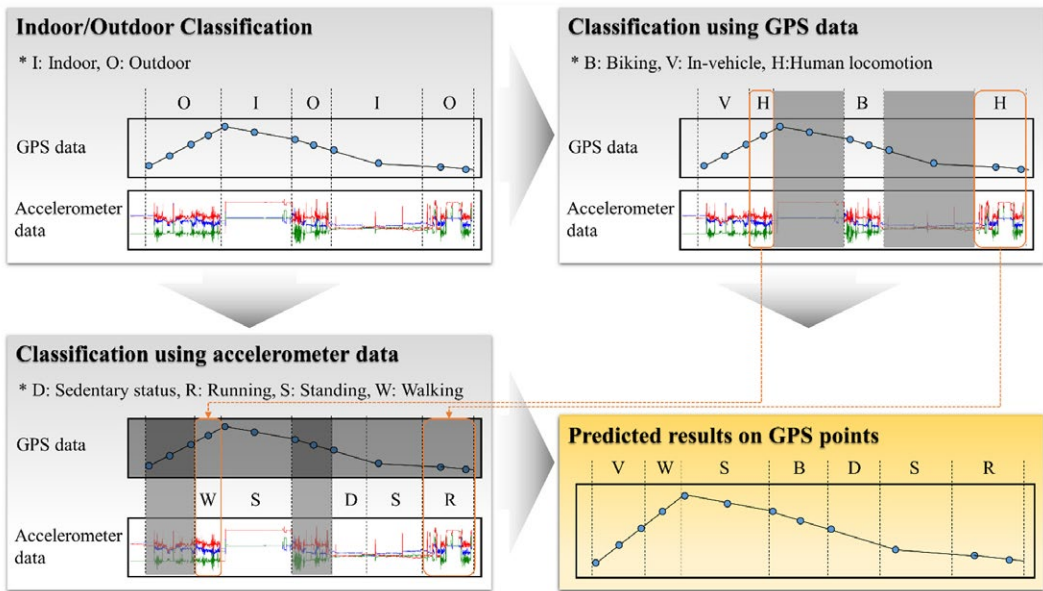
**FIGURE 2** Details of the three classification components and hierarchical classification of human locomotion

whereas accelerometer data are commonly used for classifying various types of non-vehicle-based human movement. In particular, some PA types (e.g., standing, sitting) cannot be characterized using GPS data, but can be identified using accelerometer data. A hierarchical classification approach thus allows for the identification of these PA types among non-vehicle-based movements at a lower level using accelerometer data (Figure 1). Apart from these two classification processes, GPS points are also classified into indoor or outdoor points at the top level to facilitate accurate classification processes at lower levels. For example, in an indoor environment, basic types of PA like walking and standing are likely to occur whereas in-vehicle status and biking are unlikely to take place indoors. Therefore, indoor GPS points should go through a classification process that involves the analysis of accelerometer data. In this study, indoor GPS points and accelerometer data are processed at the lower level through the "Classification using accelerometer data" process for predicting specific indoor PA types. To classify whether a setting is indoor or outdoor, we only used measures derived from the GPS data that do not come with additional sensor data (e.g., camera images, light-sensor levels) to enhance the performance of the methods implemented in previous studies (Lam et al., 2013; Tandon, Saelens, Zhou, Kerr, & Christakis, 2013).

Among the three-level classification processes shown in Figures 1 and 2, the "Indoor/Outdoor Classification" process first partitions GPS trajectories into indoor or outdoor segments, since indoor GPS segments are ineffective and unreliable for extracting any movement derivatives due to the low accuracy and loss of GPS signals in indoor environments. Outdoor GPS segments identified in this process will then go through the "Classification using GPS data" process, which classifies GPS segments into vehicle-based movement (riding in a vehicle, biking) and non-vehicle-based human movement (human locomotion). The indoor/outdoor classification and classification using GPS data processes apply the concept of GPS trajectory operators (Laube et al., 2007), as discussed in Section 3.3. On the other hand, indoor GPS segments are not used to extract features for any classification process, but used to provide time ranges within which specific indoor PA types like sedentary status and walking are predicted through "Classification using accelerometer data" (Lee & Kwan, 2017). The pre-processing is followed by classification using GPS or accelerometer data to filter out low-quality GPS points. Among the three identified classes in the "Classification using GPS data" process, PA is further classified into specific PA types such as running

and walking through the "Classification using accelerometer data" process. Through these three-level classification processes, six different PA classes are identified: traveling in a vehicle, biking, running, walking, standing, and being sedentary.

## 3.1 | Data description

Datasets from four different sources were used for the three-level classification processes (see Table 2). First, since there are no publicly available GPS data that can be used to classify indoor versus outdoor trajectories, GPS data were collected from three persons living in highly or moderately urbanized areas for seven days with horizontal dilution of precision (HDOP) and "indoor" or "outdoor" labels. Dilution of precision represents the precision of positional measurement, and specifically, HDOP indicates the geometric quality of the horizontal positions of GPS data. A large discrepancy in GPS accuracy between highly and moderately urbanized areas might affect the indoor/outdoor identification. Thus, GPS and accelerometer data were collected from one subject living in a highly urbanized area and two subjects living in a moderately urbanized area. These data were used in the test phase reported in Section 4.1. Only two days of GPS trajectories, which have many combinations of indoor and outdoor trips, from each of the three subjects were also used in the indoor/outdoor classification in Section 3.3.

Second, version 1.3 of the Geolife project GPS dataset (Microsoft Research Asia) is used for the classification using GPS data. The GPS dataset was collected from 182 subjects from April 2007 to August 2012 (Zheng, Li, Chen, Xie, & Ma, 2008; Zheng et al., 2010). Most of these GPS trajectories were tracked at short time intervals (e.g., 1 to 5 s). Only a part of the GPS data from 73 participants have labels of 11 kinds of transport modes. Among them, GPS trajectories labeled with seven transport modes—train, bus, car, taxi, biking, walking, and running—are used in the study. The total number of GPS points with these seven transport-mode labels is 5,063,475, and walking (35%) comprises a large proportion of the GPS points whereas running (0.03%) constitutes the smallest portion.

Third, due to the lack of GPS data for the running mode, open GPS data recorded when people were running were obtained from OpenStreetMap (OSM) and TrackProfiler and used in this study. Train, taxi, bus, and car transport modes are merged into one motorized mode, and walking and running (GPS data from Geolife, OSM, and TrackProfiler) are grouped into one human-locomotion class in the study.

**TABLE 2** Description of GPS and accelerometer datasets

| Dataset | No. of instances | Attributes | Mode | Sampling rate |
|---|---|---|---|---|
| GPS dataset collected from two subjects | 58,172 | Person ID, longitude, latitude, timestamp, HDOP, indoor/outdoor mode | Indoor (61%), outdoor (39%) | 1 s |
| Geolife GPS dataset v1.3 (Microsoft Research Asia) | 5,063,475 | Person ID, longitude, latitude, timestamp, transport mode | Train (5%), bus (17%), car (8%), taxi (7%), bike (28%), walk (35%), run (0.03%) | Mostly 1–5 s |
| GPS dataset from OSM and TrackProfiler | 10,435 | Longitude, latitude, timestamp | Run (100%) | Variable, 1–17 s |
| Wireless sensor data mining (WISDM)'s accelerometer dataset v1.1 | 1,098,207 | Person ID, physical activity mode, timestamp, acceleration $(x, y, z)$ | Walking (39%), running (31%), upstairs (11%), downstairs (9%), sitting (6%), standing (4%) | 20 Hz (20 samples/s) |

Finally, version 1.1 of the WIreless Sensor Data Mining (WISDM) accelerometer dataset is used for the classification using accelerometer data process (Kwapisz et al., 2011). Three-axis acceleration was collected from 36 persons with timestamp under laboratory conditions and recorded at 20 Hz. Each instance (observation) has a PA label (e.g., walking, running, or standing). Since movement on stairs accounts for a small portion of the daily PA of the subjects, movements upstairs and downstairs are merged into walking in this study.

## 3.2 | GPS trajectory operators

GPS trajectory operators (Laube et al., 2007) are used to extract features from GPS trajectories in the indoor/outdoor classification and classification using GPS data processes. Different levels of operators for GPS trajectories were proposed by introducing map algebra operations to capture dynamic movement characteristics over space and time. An operator performs specific mathematical and/or logical analysis or calculation. For instance, in map algebra, there are four types of operation—local, focal, zonal, and global—for spatial analysis to produce a resulting map using raster data (Tomlin, 1990). Among the four different levels of GPS operators (instantaneous, interval, episodal, and global), instantaneous and interval operators are used to acquire the local and focal characteristics of dynamic movement and to calculate movement derivatives (speed, acceleration) for feature extraction. In machine learning, feature extraction is a process for deriving values (features) from the data to train a model. Movement derivatives at the instantaneous level are calculated between two consecutive GPS points to capture local movement characteristics pertinent to each of these GPS points. For instance, assuming that a given GPS trajectory is recorded at a 2 s interval (as shown in Figure 3), the instantaneous velocity of the GPS point $p_3$ is calculated based on the distance and 2 s time interval ($\delta t = 2$) between $p_3$ and $p_4$.

Compared to instantaneous movement derivatives, movement derivatives at the interval level are calculated taking into account both the spatial and temporal dimensions. A distance or time window moves along a given GPS trajectory to calculate interval movement derivatives, including average velocity and maximum acceleration, considering all GPS points within the range of the distance or time window. The calculated interval derivatives particularly reflect distinguishable spatial and temporal variations in the movement characteristics among different transport modes, taking into account traffic conditions. For example, assuming that the transport mode of the given GPS trajectory in Figure 3 is "car" in smooth traffic flow, the average velocity of the GPS point $p_7$ within the 10 m distance window is calculated based on the velocity values of $p_6$, $p_7$, and $p_8$, whereas the average velocity of the same GPS point $p_7$ within the 10 s time window is calculated based on the velocity values of $p_5$, $p_6$, $p_7$, $p_8$, and $p_9$. However, if the traffic was heavy around $p_7$, the number of GPS points considered for the 10 m distance
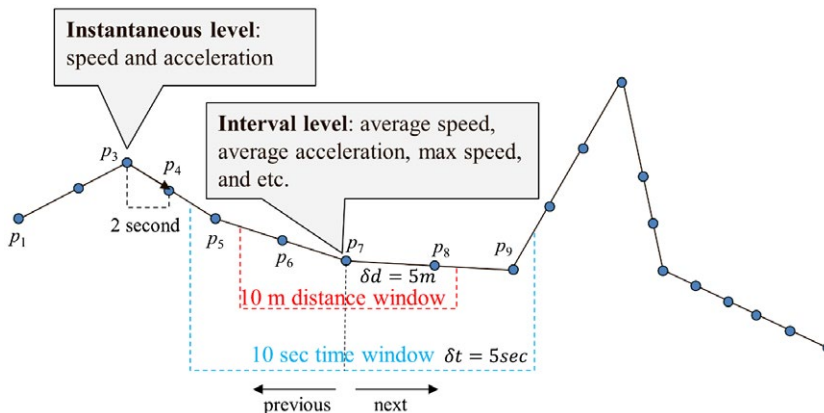


**FIGURE 3** Movement derivatives calculated at instantaneous and interval levels

window would be more than the three points due to the many stationary GPS points near $p_7$, even though the size of the window is still the same as that in a smooth traffic flow. On the other hand, the 10 s time window is likely to shrink because of many stationary GPS points around $p_7$ in such heavy traffic, which can result in a drastic drop in average velocity. In this way, interval movement derivatives using distance and time windows can capture different focal characteristics of each transport mode in a complementary manner, in order to improve performance.

## 3.3 | Indoor/outdoor classification

Due to the limited usefulness of indoor GPS data, outdoor GPS points need to be separated from indoor GPS points. Extreme gradient boosting (Chen & Guestrin, 2016), which follows the principles of gradient boosting (Friedman, 2001), is used to classify these two labels (indoor versus outdoor) based on three selected features derived from the GPS data collected from three subjects. Boosting is a forward stage-wise optimization approach that uses the votes of each weak classifier, learned at every iteration, to generate a strong classifier. Gradient boosting uses tree models as weak classifiers and generates a strong classifier based on the notion of gradients in such a way that the loss function is minimized. Extreme gradient boosting particularly takes regularization into account to control over-fitting for improved performance. Velocity, acceleration, HDOP, and the number of missing points are taken into account in feature extraction, because GPS points may have a considerable number of errors that manifest as position and velocity spikes, and may not have been properly recorded due to signal loss. Only three important features—average velocity and HDOP at the 180 s interval level, and number of missing points at the 120 s interval level—are selected among all 31 features considering both instantaneous and interval levels of the GPS trajectory operators (Table A, Appendix) for the indoor/outdoor classification.

## 3.4 | Pre-processing of indoor and outdoor GPS points

Pre-processing of GPS data involves determining valid indoor and outdoor GPS points. For the identified indoor GPS points, longitude and latitude values with high HDOP (e.g., HDOP>8 in this study) are replaced by those of previous GPS points with low HDOP so that many arbitrarily dispersed GPS spikes can be displaced on a previously tracked place. Missing indoor GPS points due to signal loss are substituted by the positions of previous GPS points with the indoor label. Regarding the identified outdoor GPS points, GPS points with high HDOP are replaced by those estimated using a Kalman filter to filter out noisy records and produce accurate GPS points. The application of the Kalman filter is described in Section 4.1.

## 3.5 | Classification using GPS data

Since the Geolife GPS data are partitioned into many segments for each person and travel mode, two or more continuous trips with less than a 1 min time gap need to be merged into one segment as an analytic unit for feature extraction. The merging process considers the transition between two different modes, like walking and riding a bus, to improve the identification of transport modes. For example, Person 1 and Person 2 traveled with different transport modes on 04/03/2008 and 08/27/2011, and different trips are divided into segments as shown in Table 3. In this case, because the first three segments of Person 1 from the first row to the third row happened continuously, they are concatenated into one segment. Further, since an important goal of the classification using the GPS data is to accurately separate traveling in a vehicle (in-vehicle mode) and biking from PA (like running) for better classification in the next phase using accelerometer data, consecutive segments of walking/running and other transport modes are merged into one concatenated segment for feature extraction.

In addition, valid GPS segments are selected using the following criteria: (1) segments that consist of GPS points with an average recorded time interval of less than 10 s; and (2) segments with a sum of recorded time of more than 3 min. GPS points recorded at longer time intervals cannot ensure consistency of the features, so the

**TABLE 3** Examples of segments in Geolife GPS data

| Row ID | Person ID | Start time | End time | Transport mode |
|--------|-----------|------------|----------|----------------|
| 1 | 1 | **04/03/2008 11:32:24** | **04/03/2008 11:46:14** | **walk** |
| 2 | 1 | **04/03/2008 11:47:14** | **04/03/2008 11:55:07** | **taxi** |
| 3 | 1 | **04/03/2008 11:55:24** | **04/03/2008 12:01:49** | **taxi** |
| 4 | 1 | 04/03/2008 16:00:00 | 04/04/2008 04:13:22 | train |
| ... | ... | ... | ... | ... |
| 101 | 2 | 08/27/2011 06:13:01 | 08/27/2011 08:01:37 | walk |
| 102 | 2 | 08/27/2011 15:01:59 | 08/27/2011 15:31:43 | walk |

*Note.* Bold fonts indicate trips considered continuous and concatenated into one segment.

segments with GPS points having such long time intervals are excluded. Further, segments with recorded time of less than 3 min are excluded due to the use of the interval-level operator with largest time window of 3 min.

At the instantaneous and interval levels, a total of 73 features are generated for each GPS point, as shown in Table 4. Two movement derivatives, velocity and acceleration, are calculated at the instantaneous level, and a set of five movement derivatives—average velocity, average acceleration, maximum velocity, maximum acceleration, and rate of change in velocity (Zheng et al., 2010)—is calculated for each time window (10 s, 20 s, and so on) and distance window (10 m, 20 m, and so on) at the interval level. Besides instantaneous and interval derivatives, the recorded hour of each GPS point is also extracted and used as a feature due to its important role in separating time-constrained activities. For instance, people are likely to engage in biking and running during the daytime, and the use of motorized-transport modes like bus and car is usually at its peak during rush hours.

Random forest (Breiman, 2001) is used to identify vehicle-based movement (traveling in a vehicle and biking) and non-vehicle-based human movement (human locomotion) based on features extracted from the Geolife GPS dataset. The mean decrease Gini index is measured for each feature (Table B, Appendix) to examine the important features. The most important movement derivatives are average velocity and maximum velocity at the interval level, whereas no movement derivatives at the instantaneous level show high importance.

## 3.6 | Classification using accelerometer data

Running, walking, standing, and sedentary status are identified in the classification using accelerometer data based on the approach developed by Lee and Kwan (2017), which has 99.03% predictive accuracy when random forest is used. The time range of indoor GPS points delineates the boundary within which such specific PA types are identified using accelerometer data. The identified human-locomotion label in the classification using GPS data also goes through the classification component using accelerometer data to generate specific PA types.

## 4 | RESULTS

### 4.1 | Performance of PA and in-vehicle status classification algorithm on real-world data

The PA and in-vehicle classification algorithm was implemented using the R statistical computing software (R Foundation for Statistical Computing, 2013). R supports machine learning model packages called "xgboost" (Chen & He, 2018) and "randomForest" (Liaw & Wiener, 2002) for extreme gradient boosting and random forest used in the indoor/outdoor classification and classification using GPS data, respectively. Since the time and distance windows at the interval level demand intensive computation in the indoor/outdoor classification and classification using GPS data processes, the "doParallel" package, which deploys parallel computing through exploiting multi-cores, was used to improve computational performance (Calaway, Weston, & Tenenbaum, 2014). R was also used to extract all the features from the GPS and accelerometer data.

**TABLE 4** All 73 features for classification using GPS data (bold font: high importance)

| GPS trajectory operator | Movement derivative | |
| --- | --- | --- |
| Instantaneous | Velocity | |
| | Acceleration | |
| | Recorded time (h) | |
| Interval | *Time window* | |
| | 10 s | **Average velocity,** average acceleration, max velocity, max acceleration, change rate of velocity |
| | 20 s | **Average velocity**, average acceleration, max velocity, max acceleration, change rate of velocity |
| | 30 s | **Average velocity**, average acceleration, max velocity, max acceleration, change rate of velocity |
| | 60 s | **Average velocity**, average acceleration, **max velocity**, max acceleration, change rate of velocity |
| | 90 s | **Average velocity**, average acceleration, **max velocity**, max acceleration, **change rate of velocity** |
| | 120 s | **Average velocity**, average acceleration, **max velocity**, max acceleration, **change rate of velocity** |
| | 180 s | **Average velocity**, average acceleration, **max velocity**, **max acceleration**, **change rate of velocity** |
| | *Distance window* | |
| | 10 m | Average velocity, average acceleration, **max velocity**, max acceleration, change rate of velocity |
| | 20 m | Average velocity, average acceleration, max velocity, max acceleration, change rate of velocity |
| | 30 m | Average velocity, average acceleration, max velocity, max acceleration, change rate of velocity |
| | 40 m | **Average velocity**, average acceleration, max velocity, max acceleration, change rate of velocity |
| | 50 m | Average velocity, average acceleration, **max velocity**, max acceleration, change rate of velocity |
| | 100 m | **Average velocity**, average acceleration, **max velocity**, max acceleration, change rate of velocity |
| | 200 m | **Average velocity**, average acceleration, **max velocity**, max acceleration, change rate of velocity |

The performance of the two learning models generated in the indoor/outdoor classification and classification using GPS data processes was evaluated using 10-fold cross-validation. For the classification using GPS data, only 218,342 instances were randomly sampled and used for the cross-validation, taking into account the original proportion of each label, due to the considerable computation time when all instances are input. The indoor/outdoor classification model using extreme gradient boosting showed 99.56% predictive accuracy in total with 500 iterations, and both indoor and outdoor GPS points were classified correctly with over 99% accuracy, as shown in Table 5.

For the classification using GPS data, random forest achieved 94.47% predictive accuracy with 500 trees. The confusion matrix in Table 6 demonstrates that human locomotion and in-vehicle classes are identified with a high accuracy of over 90%. Biking has the lowest classification accuracy (77.90%), and particularly, 1,833 and 1,061 instances with biking labels were incorrectly predicted as human locomotion and in-vehicle, respectively. Human locomotion shows the highest predictive accuracy (96.57%) among the three classes, which is, in turn, expected to deliver more accurately classified human locomotion instances to be further identified as one of the four PA types in the classification using accelerometer data.

To validate the proposed algorithm on real-world data, GPS and accelerometer data collected from three subjects in free-living conditions were tested. These three subjects are all male, young adults, and undergraduate or graduate students with no health issues. College students are appropriate subjects, who may take advantage of various opportunities in urban areas. Particularly, the three subjects often engage in MVPA in their daily lives by walking on campus, exercising, or running/biking around their residential area. They were given smartphones or used their own smartphones to record both GPS tracks and accelerometer data during a 7-, 8-, or 28-day period, as shown in Table 7. Subject 1 particularly participated in data collection for a duration of almost 1 month to explore more activities with predicted results through visualization. All three subjects were asked to collect GPS and accelerometer data, carrying the phone in their right pants' pocket. However, Subject 3 carried the provided smartphone in one of his jacket pockets by mistake, which resulted in poor predictive accuracy in the classification using accelerometer data, so the collected accelerometer data from Subject 3 were not used in this study. In addition, the accelerometer data from Subject 2, initially collected with an inexpensive phone (LG Realm) for a week, were all erroneous. All three axes had mere variations of values, and the predicted results were mostly sitting

**TABLE 5** Confusion matrix of indoor/outdoor classification using extreme gradient boosting

|  |  | Actual class | |
|---|---|---|---|
|  |  | Indoor | Outdoor |
| Predicted class | Indoor | 35,140 | 106 |
|  | Outdoor | 147 | 22,740 |
| Accuracy (%) |  | 99.58 | 99.54 |

**TABLE 6** Confusion matrix of classification component using GPS data

|  |  | Actual class | | |
|---|---|---|---|---|
|  |  | Human locomotion | In-vehicle | Biking |
| Predicted class | Human locomotion | 89,593 | 5,755 | 1,833 |
|  | In-vehicle | 2,949 | 106,466 | 1,061 |
|  | Biking | 235 | 250 | 10,200 |
| Accuracy (%) |  | 96.57 | 94.66 | 77.90 |

**TABLE 7** Description of GPS and accelerometer data collected from three subjects under free-living conditions

| Subject ID | Device | Data | Recording time |
|------------|--------|------|----------------|
| 1 | LG G3, Samsung Galaxy Alpha | GPS and accelerometer | 28 days |
| 2 | Samsung Galaxy Alpha | GPS and accelerometer | 8 days |
| 3 | LG Realm | GPS | 7 days |

status even when the subject walked. Hence, Subject 2 was asked to collect data again using a Samsung Galaxy Alpha, which led to better classification results. GPS trajectories were recorded at a 1-s interval with HDOP. In addition, activity diaries that recorded subjects' transport modes and PA types in detail were also collected. Since the developed algorithm needs to be evaluated based on very detailed records of the subjects' activities, transport modes, and PA types, each subject was asked to elaborate his activity diaries for each second by checking a visualized GPS trajectory. The elaborated 3 day activity diaries, GPS trajectories, and accelerometer data were then used as a test dataset for validating the algorithm.

Since some GPS points were excluded due to their low accuracy and inherent positional errors, the Kalman filter was used to accurately estimate the latitude and longitude values of invalid GPS points. The Kalman filter is one of the most widely used and best-performing smoothing methods for mitigating GPS random errors that influence the accuracy of derived measures from GPS points (Jun et al., 2006; Grewal et al., 2011). The function "dlmSmooth" in R for Kalman filtering was utilized in this study to estimate missing or excluded GPS points and smooth GPS trajectories (Petris, 2009). With regard to the Kalman filter, parameter values suggested by a simple dynamic linear model of Petris (2009) were applied for variance of observation noise and variance of systematic noise at default [Table C(c), Appendix]. Higher variances of systematic noise were also tested [Table C(b),(d), Appendix] to observe the effects of weak (variance of systematic noise = 1.0) or strong (variance of systematic noise = 0.01) filters, which make the test GPS trajectories less or more deviated from the original records (Figure 4). Kalman filtering with a variance of systematic noise of 0.1 had the best predictive accuracy and thus was applied to the test GPS trajectories for further evaluation of the PA and in-vehicle status classification algorithm.

The performance of the algorithm on real-world GPS trajectories achieved 96.20% accuracy in the identification of the six PA types (Table 8). The PA type identified with the lowest predictive accuracy was running (69.98%); 189 GPS points were wrongly classified as biking in the classification using GPS data. Walking and standing were the two most accurate types in the results (98.25% and 97.83%, respectively). Sensitivity, specificity, positive predictive value, and negative predictive value for the algorithm were also calculated (see Table 9).

## 4.2 | Predicted PA types and in-vehicle status and its visualization with GPS trajectories

To quantitatively validate the PA and in-vehicle status classification algorithm, PA types and in-vehicle status predicted by the developed algorithm were combined with the collected GPS points based on the timestamp in both the accelerometer data and GPS trajectories (see Figure 5). GPS trajectories coupled with the predicted PA types were visualized on a map using ArcGIS 10.4 to assess the predicted results of continuous trips and their relationships with ambient geographic contexts.

Car, bus, and biking mostly showed correct classification. Although there are waiting points with slow speeds around road intersections and bus stops, the prediction of the car and bus travel modes showed promising results. However, compared to bus in the moderately urbanized area, some trajectories of traveling by bus in the urban area were incorrectly classified as biking. Biking also showed some wrong classification results as in-vehicle, which can be explained by its moderate predictive accuracy in Section 4.1. Walking, sedentary status (sitting), and standing were also accurately classified by random forest through the classification using accelerometer data. Data on running were collected for a short time (e.g., 10 min) from one subject, and some GPS points during running were
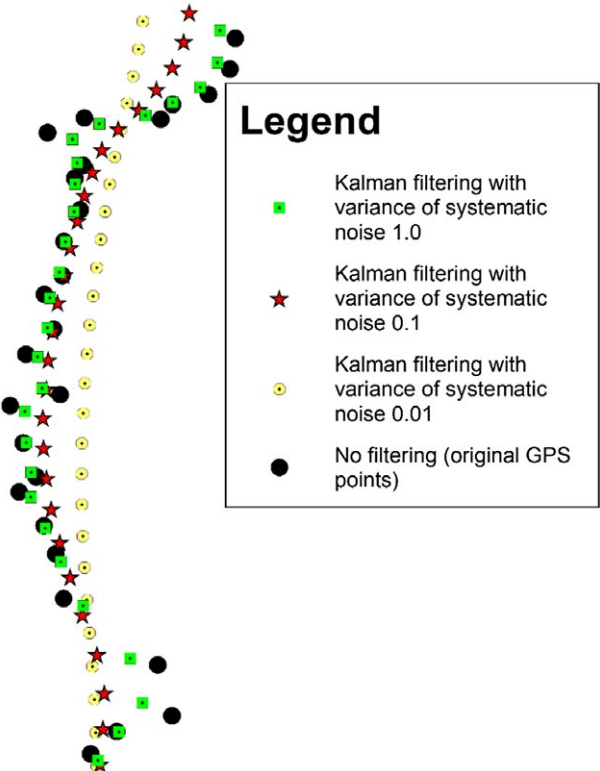
**FIGURE 4** Testing Kalman filtering on GPS trajectories with three different variance values for systematic noise

**TABLE 8** Confusion matrix of PA and in-vehicle status classification on free-living condition GPS and accelerometer dataset

|  |  | *Actual class* |  |  |  |  |  |
| --- | --- | --- | --- | --- | --- | --- | --- |
|  |  | Running | Walking | Sitting | Standing | In-vehicle | Biking |
| Predicted class | Running | 464 | 19 | 0 | 0 | 0 | 0 |
|  | Walking | 3 | 6,725 | 1,363 | 8 | 27 | 0 |
|  | Sitting | 0 | 80 | 65,889 | 0 | 101 | 2 |
|  | Standing | 0 | 0 | 1,019 | 2,889 | 17 | 0 |
|  | In-vehicle | 7 | 19 | 47 | 56 | 1,662 | 175 |
|  | Biking | 189 | 2 | 0 | 0 | 0 | 1,779 |
| Accuracy (%) |  | 69.98 | 98.25 | 96.44 | 97.83 | 91.98 | 90.95 |

wrongly classified as biking. An above-ground subway trip, which the proposed algorithm did not take into account for the in-vehicle class, was detected as in-vehicle. A series of daily activities was also examined, as shown at the bottom right of Figure 5. One subject went to a place by car for recreational purposes and for this person, walking, sitting, and standing were represented at the recreational place indoors and outdoors.

**TABLE 9** Sensitivity, specificity, positive predictive value, and negative predictive value for PA and in-vehicle status classification algorithm

| Measures | Running | Walking | Sitting | Standing | In-vehicle | Biking |
|---|---|---|---|---|---|---|
| Sensitivity (%) | 69.98 | 98.25 | 96.44 | 97.83 | 91.98 | 90.95 |
| Specificity (%) | 99.98 | 98.15 | 98.71 | 98.70 | 99.62 | 99.76 |
| Positive predictive value (%) | 96.07 | 82.76 | 99.72 | 73.61 | 84.54 | 90.31 |
| Negative predictive value (%) | 99.76 | 99.84 | 85.25 | 99.92 | 99.82 | 99.78 |



**FIGURE 5** Visualization of predicted six classified PA types and in-vehicle status in free-living conditions

## 5 | DISCUSSION AND CONCLUSIONS

The hierarchical classification algorithm developed using publicly available GPS and accelerometer data in this study showed excellent performance in accurately predicting in-vehicle status and PA types from real-world GPS and accelerometer data. Publicly available GPS and accelerometer data were shown to have potential for automatically and accurately classifying people's in-vehicle status and PA types to conduct research using objectively identified PA. The results also indicated that without any segmentation methods, in-vehicle status, biking, and PA can be successfully classified with GPS trajectory operators. Based on our brief exploration, it was also revealed that the relatively low accuracy of biking classification in Table 6 was due to some GPS points deviating from their

original tracks, presumably caused by positional errors and actual biking patterns with slow speed in business areas and possibly around destinations, even though we did not specifically analyze these in this study.

None of the instantaneous-level features were important in the classification using GPS data. This indicated that the interval-level movement descriptors calculated based on different sizes of time and distance windows outperform the instantaneous-level movement descriptors calculated from two consecutive GPS points, which have been widely used in many transport-mode classification studies. In addition, 20 important features in Table 4 suggest that not only the time dimension but also the spatial dimension should be considered to extract features. Average and maximum velocity calculated through the time and distance windows played an important role in better predicting in-vehicle status, biking, and human locomotion.

Using heterogeneous sensor datasets to classify in-vehicle status and PA with hierarchical processes is novel and has considerable potential for application in a wide range of domains. For instance, Prelipcean, Gidyfalvi, and Susilo (2017) highlighted that interdisciplinary solutions regarding travel-mode detection should not be limited to one research domain. They also emphasized that the current research trend of validating new algorithms and using datasets that cannot be shared widely hinders interdisciplinary studies. In this regard, this study provided guidance for researchers in transportation and PA research to design a classification algorithm for travel modes in an extensible manner by training models using publicly available heterogeneous sensor datasets. In addition, we will make the algorithm available to interested researchers or practitioners upon request via a webpage. It is also imperative that predicted PA types and in-vehicle status with GPS trajectories potentially enable the prediction of other activity-related or contextual information, including the activity itself and activity purposes, useful for addressing the UGCoP. For instance, the GPS points of a subject around a bus stop, coupled with predicted "standing" or "sitting," can be interpreted as waiting time for a bus. Predicted labels of "walking" can also be further separated by purpose through systematic assumption and logical reasoning if coupled GPS tracks are along certain contexts with specific purposes, like recreational facilities. Such inferred information from the predicted results, therefore, can help advance PA and transportation research. In addition, this study provides insights into how PA and transport-mode classification can be applied to various research domains involving human mobility analysis (e.g., air quality; Tainio et al., 2016), especially with regard to how the algorithm works on sensor datasets collected in free-living conditions and how the predicted results are represented.

In this study, GPS and accelerometer data collected from three subjects were classified using an algorithm based on the hierarchical classification processes. The classification results were visualized to enhance understanding of the subjects' daily activities and travel. The algorithm classified all activities with a series of plausible in-vehicle status and PA types with high accuracy: riding in a car or bus in the moderately urbanized area, biking and sitting in a building, and standing around a bus stop. Especially, the interval level of the GPS trajectory operators in the classification using GPS data helped to capture the characteristics of car and bus riding within different time and distance windows. It contributed to the correct classification of waiting points with slow speeds around road intersections and bus stops. Further, uncontrolled movement during a recreational activity was plausibly interpreted as a combination of different postures, like walking, standing, and sitting, which were characteristic of the activity. More important, the algorithm was able to detect indoor PA with high predictive accuracy.

The high performance of the developed algorithm, however, can be achieved only when the quality of collected sensor data is reliable. Low quality of built-in accelerometer sensors is likely to record erratic acceleration values, which may not be helpful for calculating features for accurate PA classification (Lee & Kwan, 2017). The placement of smartphones can also greatly affect the results of classification using accelerometer data. Therefore, a smartphone screening procedure and instruction for participants on the placement of smartphones will be necessary to assure the performance of the algorithm. Additionally, none of the publicly accessible data used to train the models includes any personal information, like demographic information or socioeconomic status, which indicates the specific population subgroups for whom the algorithm works best. Personal-level characteristics can account for internal variations among individuals in terms of not only age, gender, and race, but also daily PA patterns, health conditions, and transport modes, which in turn collectively constitutes population characteristics of

the training data. However, since the Geolife GPS dataset, OSM GPS dataset, and WISDM accelerometer data do not provide any individual or population characteristics, whether the algorithm can achieve the same performance for other population groups is uncertain.

The classification algorithm developed in this study has some limitations that need to be addressed in future research. First, additional processing should be applied in order to minimize incorrectly classified PA types. For instance, bus riding in urban areas has inconsistent classification results, whereas bus riding in moderately urbanized areas was perfectly classified. The incorrect classifications were largely the result of the low positional accuracy of the GPS points. In addition, a few wrongly predicted PA types also appeared intermittently in the middle of a trip. In this case, post-processing is needed to replace the erroneous PA labels with the correct ones.

In addition, the design of the classification algorithm needs to be improved for higher predictive accuracy. The hierarchical classification propagated an error downwards onto other classification processes. For example, in the PA prediction using the accelerometer data, some points of running were previously misclassified as biking in the classification using GPS data; thus, the PA classification using accelerometer data could not intervene in the prediction and change the biking label to the running label, even though the PA classification using accelerometer data correctly classfied it. Thus, future research should address how effectively the predicted results from the hierarchical classification processes can be exploited.

Finally, the motorized-transport mode should be classified further to account for different motorized modes. This study combined train, bus, car, and taxi into a single motorized-transport mode because the subdivision of the motorized transport does not have any significant meaning in terms of intensity and type of PA. However, for other research domains such as transport or social science, the differentiation of motorized-transport modes, whether public or private, can be important (Biljecki et al., 2013; Feng & Timmermans, 2013; Shafique & Hato, 2016). Thus, the single motorized-transport mode in this study needs to be subdivided and specified for this kind of research in the future.

**ORCID**

*Kangjae Lee* iD http://orcid.org/0000-0002-2857-6496
*Mei-Po Kwan* iD http://orcid.org/0000-0001-8602-9258

**REFERENCES**

Almanza, E., Jerrett, M., Dunton, G., Seto, E., & Pentz, M. A. (2012). A study of community design, greenness, and physical activity in children using satellite, GPS and accelerometer data. *Health & Place*, *18*, 46–54.

Anguita, D., Ghio, A., Oneto, L., Parra, X., & Reyes-Ortiz, J. L. (2012). Human activity recognition on smartphones using a multiclass hardware-friendly support vector machine. In *Proceedings of the Fourth International Workshop on Ambient Assisted Living* (pp. 216–223). Vitoria-Gasteiz, Spain.

Arif, M., Bilal, M., Kattan, A., & Ahamed, S. I. (2014). Better physical activity classification using smartphone acceleration sensor. *Journal of Medical Systems*, *8*, 95.

Biljecki, F., Ledoux, H., & Van Oosterom, P. (2013). Transportation mode-based segmentation and classification of movement trajectories. *International Journal of Geographical Information Science*, *27*, 385–407.

Bolbol, A., Cheng, T., Tsapakis, I., & Haworth, J. (2012). Inferring hybrid transportation modes from sparse GPS data using a moving window SVM classification. *Computers, Environment & Urban Systems*, *36*, 526–537.

Boruff, B. J., Nathan, A., & Nijnstein, S. (2012). Using GPS technology to (re)-examine operational definitions of 'neighbourhood' in place-based health research. *International Journal of Health Geographics*, *11*, 22.

Breiman, L. (2001). Random forests. *Machine Learning*, *45*, 5–32.

Broach, J., Dill, J., & Gliebe, J. (2012). Where do cyclists ride? A route choice model developed with revealed preference GPS data. *Transportation Research Part A: Policy & Practice*, *46*, 1730–1740.

Browning, M., & Lee, K. (2017). Within what distance does "greenness" best predict physical health? A systematic review of articles with GIS buffer analyses across the lifespan. *International Journal of Environmental Research & Public Health*, *14*, 675.

Calaway, R., Weston, S., & Tenenbaum, D. (2014). *Do Parallel: Foreach parallel adaptor for the parallel package*. Retrieved from https://rdrr.io/cran/doParallel/

Chen, T., & Guestrin, C. (2016). *Xgboost: A scalable tree boosting system*. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (pp. 785–794). San Francisco, CA: ACM.

Chen, T., & He, T. (2018). *Xgboost: eXtreme gradient boosting*. Retrieved from https://cran.r-project.org/web/packages/xgboost/vignettes/xgboost.pdf

Cooper, A. R., Page, A. S., Wheeler, B. W., Griew, P., Davis, L., Hillsdon, M., & Jago, R. (2010). Mapping the walk to school using accelerometry combined with a global positioning system. *American Journal of Preventive Medicine*, *38*, 178–183.

Council on Tall Buildings and Urban Habitat. (2018). *The Skyscraper Center*. Retrieved from https://www.skyscrapercenter.com/city/chicago

Diez Roux, A. V., & Mair, C. (2010). Neighborhoods and health. *Annals of the New York Academy of Sciences*, *1186*, 125–145.

Dodge, S., Bohrer, G., Weinzierl, R., Davidson, S. C., Kays, R., Douglas, D., … Wikelski, M. (2013). The environmental-data automated track annotation (EnvDATA) system: Linking animal tracks with environmental data. *Movement Ecology*, *1*, 3.

Dodge, S., Weibel, R., & Forootan, E. (2009). Revealing the physics of movement: Comparing the similarity of movement characteristics of different types of moving objects. *Computers, Environment & Urban Systems*, *33*, 419–434.

Downs, J. A., & Horner, M. W. (2012). Probabilistic potential path trees for visualizing and analyzing vehicle tracking data. *Journal of Transport Geography*, *23*, 72–80.

Dumais, S., & Chen, H. (2000). Hierarchical classification of web content. In *Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* (pp. 256–263). Athens, Greece: ACM.

Eagle, N., Pentland, A. S., & Lazer, D. (2009). Inferring friendship network structure by using mobile phone data. *Proceedings of the National Academy of Sciences of the USA*, *106*, 15274–15278.

Ellis, K., Godbole, S., Marshall, S., Lanckriet, G., Staudenmayer, J., & Kerr, J. (2014). Identifying active travel behaviors in challenging environments using GPS, accelerometers, and machine learning algorithms. *Public Health*, *2*, 39–46.

Evenson, K. R., & Furberg, R. D. (2017). Moves app: A digital diary to track physical activity and location. *British Journal of Sports Medicine*, *51*, 1169–1170.

Fang, S.-H., Liao, H.-H., Fei, Y.-X., Chen, K.-H., Huang, J.-W., Lu, Y.-D., & Tsao, Y. (2016). Transportation modes classification using sensors on smartphones. *Sensors*, *16*, 1324.

Feng, T., & Timmermans, H. J. (2013). Transportation mode recognition using GPS and accelerometer data. *Transportation Research Part C: Emerging Technologies*, *37*, 118–130.

Ferreira, N., Poco, J., Vo, H. T., Freire, J., & Silva, C. T. (2013). Visual exploration of big spatiotemporal urban data: A study of New York City taxi trips. *IEEE Transactions on Visualization & Computer Graphics*, *19*, 2149–2158.

Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, *29*, 1189–1232.

Gong, H., Chen, C., Bialostozky, E., & Lawson, C. T. (2012). A GPS/GIS method for travel mode detection in New York City. *Computers, Environment & Urban Systems*, *36*, 131–139.

Gonzalez, M. C., Hidalgo, C. A., & Barabasi, A.-L. (2008). Understanding individual human mobility patterns. *Nature*, *453*, 779–782.

Gordon-Larsen, P., Nelson, M. C., & Beam, K. (2005). Associations among active transportation, physical activity, and weight status in young adults. *Obesity Research*, *13*, 868–875.

Grewal, M. S. (2011). Kalman filtering. In M. Lovric (Ed.), *International encyclopedia of statistical science* (pp. 705–708). Berlin, Germany: Springer.

Jankowska, M. M., Schipperijn, J., & Kerr, J. (2015). A framework for using GPS data in physical activity and sedentary behavior studies. *Exercise & Sport Sciences Reviews*, *43*, 48.

Jansen, M., Ettema, D., Pierik, F., & Dijst, M. (2016). Sports facilities, shopping centers or homes: What locations are important for adults' physical activity? A cross-sectional study. *International Journal of Environmental Research & Public Health*, *13*, 287.

Jun, J., Guensler, R., & Ogle, J. (2006). Smoothing methods to minimize impact of global positioning system random error on travel distance, speed, and acceleration profile estimates. *Transportation Research Record*, *1972*, 141–150.

Kwan, M.-P. (2004). GIS methods in time-geographic research: Geocomputation and geovisualization of human activity patterns. *Geografiska Annaler B*, *86*, 267–280.

Kwan, M.-P. (2012a). How GIS can help address the uncertain geographic context problem in social science research. *Annals of GIS*, *18*, 245–255.

Kwan, M.-P. (2012b). The uncertain geographic context problem. *Annals of the Association of American Geographers*, *102*, 958–968.

Kwan, M.-P. (2013). Beyond space (as we knew it): Toward temporally integrated geographies of segregation, health, and accessibility. *Annals of the Association of American Geographers*, *103*, 1078–1086.

Kwapisz, J. R., Weiss, G. M., & Moore, S. A. (2011). Activity recognition using cell phone accelerometers. *ACM SigKDD Explorations Newsletter*, *12*, 74–82.

Lachowycz, K., Jones, A. P., Page, A. S., Wheeler, B. W., & Cooper, A. R. (2012). What can global positioning systems tell us about the contribution of different types of urban greenspace to children's physical activity? *Health & Place*, *18*, 586–594.

Lam, M. S., Godbole, S., Chen, J., Oliver, M., Badland, H., Marshall, S. J., … Kerr, J. (2013). Measuring time spent outdoors using a wearable camera and GPS. In *Proceedings of the Fourth International SenseCam & Pervasive Imaging Conference* (pp. 1–7). San Diego, CA.

Laube, P., Dennis, T., Forer, P., & Walker, M. (2007). Movement beyond the snapshot-dynamic analysis of geospatial life-lines. *Computers, Environment & Urban Systems*, *31*, 481–501.

Lee, K., & Kwan, M.-P. (2017). Physical activity classification in free-living conditions using smartphone accelerometer data and exploration of predicted result. *Computers, Environment & Urban Systems*, *67*, 124–131.

Liaw, A., & Wiener, M. (2002). Classification and regression by random forest. *R News*, *2*, 18–22.

Loukaitou-Sideris, A. (2006). Is it safe to walk? Neighborhood safety and security considerations and their effects on walking. *CPL Bibliography*, *20*, 219–232.

McNamara, D. S., Crossley, S. A., Roscoe, R. D., Allen, L. K., & Dai, J. (2015). A hierarchical classification approach to automated essay scoring. *Assessing Writing*, *23*, 35–59.

Moiseeva, A., Jessurun, J., & Timmermans, H. (2010). Semi-automatic imputation of activity travel diaries: Use of global positioning system traces, prompted recall, and context-sensitive learning algorithms. *Transportation Research Record*, *2183*, 60–68.

Outdoor Foundation. (2016). *Outdoor recreation participation topline report 2016*. Retrieved from https://www.outdoor-foundation.org/pdf/ResearchParticipation2016Topline.pdf

Papinski, D., Scott, D. M., & Doherty, S. T. (2009). Exploring the route choice decision-making process: A comparison of planned and observed routes obtained using person-based GPS. *Transportation Research Part F: Traffic Psychology & Behaviour*, *12*, 347–358.

Patrick, K., Kerr, J., Norman, G., Ryan, S., Sallis, J., Krueger, I., … Ainsworth, B. (2008). Geospatial measurement & analysis of physical activity: Physical activity location measurement system (PALMS). *Epidemiology*, *19*, S63.

Perchoux, C., Chaix, B., Cummins, S., & Kestens, Y. (2013). Conceptualization and measurement of environmental exposure in epidemiology: Accounting for activity space related to daily mobility. *Health & Place*, *21*, 86–93.

Petris, G. (2009). *dlm: An R package for Bayesian analysis of dynamic linear models*. Retrieved from ftp://ftp.math.ethz.ch/sfs/pub/Software/RCRAN/web/packages/dlm/vignettes/dlm.pdf

Prelipcean, A. C., Gidyfalvi, G., & Susilo, Y. O. (2017). Transportation mode detection: An in-depth review of applicability and reliability. *Transport Reviews*, *37*, 442–464.

R Foundation for Statistical Computing. (2013). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.

Rodriguez, D. A., Cho, G.-H., Evenson, K. R., Conway, T. L., Cohen, D., Ghosh-Dastidar, B., … Lytle, L. A. (2012). Out and about: Association of the built environment with physical activity behaviors of adolescent females. *Health & Place*, *18*, 55–62.

Sallis, J., Bauman, A., & Pratt, M. (1998). Environmental and policy interventions to promote physical activity. *American Journal of Preventive Medicine*, *15*, 379–397.

Schrank, D., Eisele, B., Lomax, T., & Bak, J. (2015). *2015 urban mobility scorecard*. Retrieved from https://static.tti.tamu.edu/tti.tamu.edu/documents/mobility-scorecard-2015.pdf

Shafique, M. A., & Hato, E. (2016). Travel mode detection with varying smartphone data collection frequencies. *Sensors*, *16*, 716.

Shoval, N., & Isaacson, M. (2007). Tracking tourists in the digital age. *Annals of Tourism Research*, *34*, 141–159.

Shoval, N., Wahl, H.-W., Auslander, G., Isaacson, M., Oswald, F., Edry, T., … Heinik, J. (2011). Use of the global positioning system to measure the out-of-home mobility of older adults with differing cognitive functioning. *Aging & Society*, *31*, 849–869.

Tainio, M., de Nazelle, A. J., Gu̧tschi, T., Kahlmeier, S., Rojas-Rueda, D., Nieuwenhuijsen, M., … Woodcock, J. (2016). Can air pollution negate the health benefits of cycling and walking? *Preventive Medicine*, *87*, 233–236.

Tandon, P. S., Saelens, B. E., Zhou, C., Kerr, J., & Christakis, D. A. (2013). Indoor versus outdoor time in preschoolers at child care. *American Journal of Preventive Medicine*, *44*, 85–88.

Tomlin, C. D. (1990). *Geographic information systems and cartographic modeling*. Englewood Cliffs, NJ: Prentice Hall.

Troped, P. J., Wilson, J. S., Matthews, C. E., Cromley, E. K., & Melly, S. J. (2010). The built environment and location-based physical activity. *American Journal of Preventive Medicine*, *38*, 429–438.

van Dijk, J. (2018). Identifying activity-travel points from GPS-data with multiple moving windows. *Computers, Environment & Urban Systems*, *70*, 84–101.

Voss, C., Sims-Gould, J., Ashe, M. C., McKay, H. A., Pugh, C., & Winters, M. (2016). Public transit use and physical activity in community-dwelling older adults: Combining GPS and accelerometry to assess transportation-related physical activity. *Journal of Transport & Health*, *3*, 191–199.

Wang, J., Kwan, M.-P., & Chai, Y. (2018). An innovative context-based crystal-growth activity space method for environmental exposure assessment: A study using GIS and GPS trajectory data collected in Chicago. *International Journal of Environmental Research & Public Health*, *15*, 703.

Weiss, G. M., Lockhart, J. W., Pulickal, T. T., McHugh, P. T., Ronan, I. H., & Timko, J. L. (2016). Actitracker: A smartphone-based activity recognition system for improving health and well-being. In *Proceedings of the 2016 IEEE International Conference on Data Science & Advanced Analytics* (pp. 682–688). Montréal, Canada: IEEE.

Witayangkurn, A., Horanont, T., Ono, N., Sekimoto, Y., & Shibasaki, R. (2013). Trip reconstruction and transportation mode extraction on low data rate GPS data from mobile phone. In *Proceedings of the 13th International Conference on Computers in Urban Planning & Urban Management*. Utrecht, The Netherlands.

Xiao, Z., Wang, Y., Fu, K., & Wu, F. (2017). Identifying different transportation modes from trajectory data using tree-based ensemble classifiers. *ISPRS International Journal of GeoInformation*, *6*, 57.

Zhang, S., McCullagh, P., Nugent, C., & Zheng, H. (2010). Activity monitoring using a smart phone's accelerometer with hierarchical classification. In *Proceedings of the 2010 Sixth International Conference on Intelligent Environments* (pp. 158–163). Kuala Lumpur, Malaysia.

Zheng, Y., Chen, Y., Li, Q., Xie, X., & Ma, W.-Y. (2010). Understanding transportation modes based on GPS data for web applications. *ACM Transactions on the Web*, *4*, 1.

Zheng, Y., Li, Q., Chen, Y., Xie, X., & Ma, W.-Y. (2008). Understanding mobility based on GPS data. In *Proceedings of the 10th International Conference on Ubiquitous Computing* (pp. 312–321). Seoul, South Korea: ACM.

Zhou, X. (2014). *Investigating the association between the built environment and active travel of young adults using location based technology* (Unpublished Ph.D. Dissertation). University of Illinois at Urbana-Champaign, Urbana-Champaign, Illinois.

Zhu, Q., Zhu, M., Li, M., Fu, M., Huang, Z., Gan, Q., & Zhou, Z. (2016). Identifying transportation modes from raw GPS data. In *Proceedings of the Second International Conference of Young Computer Scientists, Engineers & Educators* (pp. 395–409). Harbin, China.

**APPENDIX**

**TABLE A**  All 31 GPS features and relative importance for indoor/outdoor classification (bold font: three selected important features)

| GPS trajectory operator | Movement derivative | | Relative importance | Importance order |
|---|---|---|---|---|
| Instantaneous | Velocity | | 0.001 | 25 |
| | Acceleration | | 0.000 | 31 |
| | HDOP | | 0.001 | 24 |
| Interval | 10 s | average velocity | 0.003 | 14 |
| | | average acceleration | 0.000 | 30 |
| | | number of missing GPS points | 0.003 | 16 |
| | | HDOP | 0.010 | 5 |
| | 20 s | average velocity | 0.002 | 20 |
| | | average acceleration | 0.000 | 29 |
| | | number of missing GPS points | 0.001 | 28 |
| | | HDOP | 0.004 | 13 |
| | 30 s | average velocity | 0.002 | 17 |
| | | average acceleration | 0.002 | 21 |
| | | number of missing GPS points | 0.009 | 6 |
| | | HDOP | 0.002 | 18 |
| | 60 s | average velocity | 0.001 | 23 |
| | | average acceleration | 0.001 | 27 |
| | | number of missing GPS points | 0.004 | 10 |
| | | HDOP | 0.003 | 15 |
| | 90 s | average velocity | 0.004 | 11 |
| | | average acceleration | 0.002 | 19 |
| | | number of missing GPS points | 0.001 | 26 |
| | | HDOP | 0.008 | 8 |
| | 120 s | average velocity | 0.020 | 4 |
| | | average acceleration | 0.002 | 22 |
| | | **number of missing GPS points** | **0.050** | **3** |
| | | HDOP | 0.006 | 9 |
| | 180 s | **average velocity** | **0.488** | **1** |
| | | average acceleration | 0.004 | 12 |
| | | number of missing GPS points | 0.009 | 7 |
| | | **HDOP** | **0.356** | **2** |

**TABLE B** All 73 GPS features and relative importance for classification using GPS data (bold font: top 20 important features)

| GPS trajectory operator | Movement derivative | | Relative importance | Importance order |
|---|---|---|---|---|
| Instantaneous | Velocity | | 1,050 | 26 |
| | Acceleration | | 300 | 70 |
| | Recorded time (h) | | 825 | 35 |
| Interval | *Time window* | | | |
| | **10 s** | **average velocity** | **776** | **37** |
| | | average acceleration | 343 | 67 |
| | | max velocity | 754 | 38 |
| | | max acceleration | 409 | 60 |
| | | change rate of velocity | 207 | 73 |
| | **20 s** | **average velocity** | **1,671** | **15** |
| | | average acceleration | 415 | 59 |
| | | max velocity | 946 | 31 |
| | | max acceleration | 456 | 55 |
| | | change rate of velocity | 297 | 71 |
| | **30 s** | **average velocity** | **1,612** | **16** |
| | | average acceleration | 440 | 57 |
| | | max velocity | 1,118 | 23 |
| | | max acceleration | 533 | 49 |
| | | change rate of velocity | 485 | 52 |

**TABLE B** (Continued)

| GPS trajectory operator | Movement derivative | | Relative importance | Importance order |
|---|---|---|---|---|
| 60 s | **average velocity** | | **4,524** | **8** |
| | average acceleration | | 524 | 50 |
| | **max velocity** | | **2,975** | **10** |
| | max acceleration | | 791 | 36 |
| | change rate of velocity | | 595 | 44 |
| 90 s | **average velocity** | | **4,641** | **7** |
| | average acceleration | | 588 | 46 |
| | **max velocity** | | **7,693** | **3** |
| | max acceleration | | 1,005 | 29 |
| | change rate of velocity | | 1,015 | 28 |
| 120 s | **average velocity** | | **9,007** | **2** |
| | average acceleration | | 598 | 43 |
| | **max velocity** | | **5,543** | **4** |
| | max acceleration | | 1,089 | 25 |
| | **change rate of velocity** | | **1,239** | **19** |
| 180 s | **average velocity** | | **9,088** | **1** |
| | average acceleration | | 737 | 40 |
| | **max velocity** | | **5,429** | **5** |
| | **max acceleration** | | **1,247** | **18** |
| | **change rate of velocity** | | **2,562** | **11** |

**TABLE B** (Continued)

| GPS trajectory operator | Movement derivative | Relative importance | Importance order |
|---|---|---|---|
| | *Distance window* | | |
| 10 m | average velocity | 1,215 | 21 |
| | average acceleration | 311 | 69 |
| | **max velocity** | **1,229** | **20** |
| | max acceleration | 389 | 62 |
| | change rate of velocity | 232 | 72 |
| 20 m | average velocity | 1,103 | 24 |
| | average acceleration | 346 | 66 |
| | max velocity | 472 | 53 |
| | max acceleration | 421 | 58 |
| | change rate of velocity | 326 | 68 |
| 30 m | average velocity | 890 | 32 |
| | average acceleration | 358 | 65 |
| | max velocity | 662 | 42 |
| | max acceleration | 460 | 54 |
| | change rate of velocity | 384 | 63 |
| 40 m | **average velocity** | **1,346** | **17** |
| | average acceleration | 370 | 64 |
| | max velocity | 841 | 34 |
| | max acceleration | 499 | 51 |
| | change rate of velocity | 448 | 56 |

**TABLE B** (Continued)

| GPS trajectory operator | Movement derivative | | Relative importance | Importance order |
| --- | --- | --- | --- | --- |
| | 50 m | average velocity | 1,209 | 22 |
| | | average acceleration | 392 | 61 |
| | | **max velocity** | **2,260** | **12** |
| | | max acceleration | 567 | 47 |
| | | change rate of velocity | 595 | 44 |
| | 100 m | **average velocity** | **1,767** | **14** |
| | | average acceleration | 556 | 48 |
| | | **max velocity** | **2,208** | **13** |
| | | max acceleration | 707 | 41 |
| | | change rate of velocity | 740 | 39 |
| | 200 m | **average velocity** | **4,840** | **6** |
| | | average acceleration | 845 | 33 |
| | | **max velocity** | **3,989** | **9** |
| | | max acceleration | 1,040 | 27 |
| | | change rate of velocity | 966 | 30 |

**TABLE C** Test results of the application of Kalman filtering to improve performance of physical activity and in-vehicle status classification

**(a) No filtering methods (predictive accuracy: 95.26%)**

| | | Actual class | | | | | |
|---|---|---|---|---|---|---|---|
| | | Running | Walking | Sitting | Standing | In-vehicle | Biking |
| Predicted class | Running | 414 | 21 | 0 | 0 | 0 | 0 |
| | Walking | 3 | 6,716 | 1,363 | 8 | 33 | 0 |
| | Sitting | 0 | 80 | 65,855 | 0 | 86 | 4 |
| | Standing | 0 | 0 | 1,019 | 2,789 | 4 | 0 |
| | In-vehicle | 15 | 28 | 81 | 156 | 1,684 | 774 |
| | Biking | 231 | 0 | 0 | 0 | 0 | 1,178 |
| Accuracy (%) | | 62.14 | 98.12 | 96.39 | 94.41 | 93.19 | 60.23 |

**(b) Kalman filtering with variance of observation noise 0.8 and variance of systematic noise 0.01 (predictive accuracy: 96.06)**

| | | Actual class | | | | | |
|---|---|---|---|---|---|---|---|
| | | Running | Walking | Sitting | Standing | In-vehicle | Biking |
| Predicted class | Running | 625 | 21 | 0 | 0 | 0 | 0 |
| | Walking | 25 | 6,708 | 1,363 | 8 | 41 | 0 |
| | Sitting | 7 | 80 | 65,924 | 0 | 390 | 4 |
| | Standing | 0 | 0 | 1,019 | 2,945 | 17 | 0 |
| | In-vehicle | 6 | 36 | 4 | 0 | 1,179 | 43 |
| | Biking | 0 | 0 | 8 | 0 | 180 | 1,909 |
| Accuracy (%) | | 94.27 | 98.00 | 96.50 | 99.73 | 65.25 | 97.60 |

**TABLE C** (Continued)

(c) Kalman filtering with variance of observation noise 0.8 and variance of systematic noise 0.1 (predictive accuracy: 96.20%)

| Predicted class | | Actual class | | | | | |
|---|---|---|---|---|---|---|---|
| | | Running | Walking | Sitting | Standing | In-vehicle | Biking |
| | Running | 464 | 19 | 0 | 0 | 0 | 0 |
| | Walking | 3 | 6,725 | 1,363 | 8 | 27 | 0 |
| | Sitting | 0 | 80 | 65,889 | 0 | 101 | 2 |
| | Standing | 0 | 0 | 1,019 | 2,889 | 17 | 0 |
| | In-vehicle | 7 | 19 | 47 | 56 | 1,662 | 175 |
| | Biking | 189 | 2 | 0 | 0 | 0 | 1,779 |
| Accuracy (%) | | 69.98 | 98.25 | 96.44 | 97.83 | 91.98 | 90.95 |

(d) Kalman filtering with variance of observation noise 0.8 and variance of systematic noise 1.0 (predictive accuracy: 95.61%)

| Predicted class | | Actual class | | | | | |
|---|---|---|---|---|---|---|---|
| | | Running | Walking | Sitting | Standing | In-vehicle | Biking |
| | Running | 331 | 8 | 0 | 0 | 0 | 0 |
| | Walking | 1 | 6,702 | 1,363 | 8 | 34 | 0 |
| | Sitting | 0 | 80 | 65,879 | 0 | 101 | 2 |
| | Standing | 0 | 0 | 1,019 | 2,790 | 5 | 0 |
| | In-vehicle | 5 | 42 | 57 | 155 | 1,667 | 401 |
| | Biking | 326 | 13 | 0 | 0 | 0 | 1,553 |
| Accuracy (%) | | 49.92 | 97.91 | 96.43 | 94.48 | 92.25 | 79.40 |