



Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

**ScienceDirect**

Procedia Computer Science 46 (2015) 1181 – 1187

**Procedia**  
Computer Science

International Conference on Information and Communication Technologies (ICICT 2014)

## Smart Phone Based Data Mining For Human Activity Recognition

Girija Chetty<sup>a,\*</sup>, Matthew White<sup>b</sup>, Farnaz Akther<sup>a</sup>

<sup>a</sup> Faculty of ESTeM, University of Canberra, Australia

<sup>b</sup> Infinity Imaging Pty. Ltd. Melbourne, Australia.

---

### Abstract

Automatic activity recognition systems aim to capture the state of the user and its environment by exploiting heterogeneous sensors, and permit continuous monitoring of numerous physiological signals, where these sensors are attached to the subject's body. This can be immensely useful in healthcare applications, for automatic and intelligent daily activity monitoring for elderly people. In this paper, we present novel data analytic scheme for intelligent Human Activity Recognition (AR) using smartphone inertial sensors based on information theory based feature ranking algorithm and classifiers based on random forests, ensemble learning and lazy learning. Extensive experiments with a publicly available database<sup>1</sup> of human activity with smart phone inertial sensors show that the proposed approach can indeed lead to development of intelligent and automatic real time human activity monitoring for eHealth application scenarios for elderly, disabled and people with special needs.

© 2015 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Peer-review under responsibility of organizing committee of the International Conference on Information and Communication Technologies (ICICT 2014)

*Keywords:* smart phone; activity recognition; machine learning; assisted living

---

### 1. Introduction

The first commercial hand-held mobile phones appeared in 1979, and since then there has been an unprecedented growth in the adoption of mobile phone technology, reaching to more than 80% of the world population by 2011<sup>2</sup>. Lately, smartphones, which are a new generation of mobile phones, are equipped with many powerful features

---

\*Corresponding author. Tel . 0412310060  
Email [girija.chetty@canberra.edu.au](mailto:girija.chetty@canberra.edu.au)

including multitasking and a variety of sensors, in addition to the basic telephony. The integration of these mobile devices in our daily life is growing rapidly, and it is envisaged that such devices can seamlessly monitor and keep track of our activities, learn from them and assist us in making decisions. Such assistive technologies can be of immense use for remote health care, for the elderly, the disabled and those with special needs, if there are autonomous and intelligent. However, currently, though there is good capacity for collecting the data with such smart devices, there is limited capability in terms of automatic decision support capability and making sense out of this large data repository. There is an urgent need for new data mining and machine learning techniques to be developed to this end. In this paper we propose a new scheme for human activity recognition using smart phone data, with potential applications in automatic assisted living technologies. Activity recognition systems aim to identify the actions carried out by a human, from the data collected from the sensors and the surrounding environment. The current smart phones have motion, acceleration or inertial sensors, and by exploiting the information retrieved from these sensors, recognition of activities and events can be recognized. Automatic recognition of activities and events is possible by processing this sensor data with appropriate machine learning and data mining approaches. Rest of the paper is organized as follows. The details of the publicly available activity recognition data set used in this work are described in Section 2. Section 3 discusses the relevant background work done in this area, and the proposed automatic activity recognition approach is discussed in Section 4. The experimental validation of the proposed approach is described in Section 5, and the paper concludes with some outcomes achieved from this work and the plan for future research.

## 2. Activity Recognition Database

For experimental validation of our approach, we used a publicly available activity recognition (AR) database<sup>1</sup>. This database includes labelled data collected from 30 subjects in age group of 19 – 48 years. Each person performed different activities wearing a smart phone around the waist, and engaged in six different activities—walking on flat ground and up and down stairs, sitting, standing, and lying down. A Samsung Galaxy S2 smartphone was used for data collection, which contains an accelerometer and a gyroscope for measuring 3-axial linear acceleration and angular velocity respectively at a constant rate of 50Hz, which is sufficient for capturing human body motion. The database consists of two data sets one raw pre-processing data from the sensors and another data set with features extracted.



Fig. 1. Block Schematic for Activity Recognition Processing

The raw data was then pre-processed by applying noise filters and then sampled with fixed-width sliding windows of 2.56 sec and 50% overlap. From each window, a vector of 17 features is obtained by calculating variables from the accelerometer signals in the time and frequency domain (e.g. mean, standard deviation, signal magnitude area, entropy, signal-pair correlation, etc.).

The other dataset consists of vectors that each contain 561 features and represent 2.56 seconds of time. Each vector encodes characteristics such as the tri axial average, maximum, and minimum acceleration and angular velocity over the given interval, as well as more complex properties such as the Fourier transform and autoregressive coefficients. We used this dataset for experimental validation of our activity recognition approach. The block schematic for the processing the data from this dataset is shown in Figure 1. Next Section discusses the related background work.

### 3. Background

The role of smartphones for automatic activity recognition can have several advantages due to easy device portability, and no requirement for additional fixed equipment that could be obtrusive and uncomfortable to the user. Other established activity recognition approaches are based on special purpose hardware set ups and body sensor networks as in <sup>5,6</sup>. It is unrealistic, however, to expect in general home settings for people to wear them for their daily activities, because of their difficulty, time and convenience to wear them on daily basis, though, such elaborate setups can enhance the activity recognition performance. Smart phones have an advantage because of their ease and convenience, along with the capability of multiple sensors on the phone, which can be exploited for activity recognition. Appropriate machine learning and data mining methods need to be developed for processing these multiple sensor signals from smartphones for automatic and intelligent activity recognition. Though there have been several machine learning methods available<sup>7, 8, 9, 10</sup>, it is not clear, which algorithm can performs better for activity recognition with smartphones. If automatic activity recognition systems can be built based on intelligent processing of multiple sensor features on smart phones, it will be a great contribution to eHealth area, particularly for remote activity monitoring and recognition in aged care and disability care sector.

In this article, we examine several new machine learning and data mining approaches based on decision trees and ensemble learning techniques including random forests and random committee, and compare them with traditional naïve Bayes classifier and unsupervised k-Means clustering approaches for processing smartphone sensor signals for activity recognition. The experimental evaluation of the proposed schemes with a publicly available smartphone activity recognition database <sup>1</sup> shows a significant improvement in recognition performance of proposed machine learning and data mining approaches, as compared to other methods proposed in the literature for smartphone based activity recognition. Next Section describes the details of machine learning techniques used for developing smartphone based automatic activity recognition system.

### 4. Proposed Data Mining Scheme

In this Section, we propose the data mining approach for classifying different activities in this work. As the dimensionality of features is very high (561 features), which can severely affect the implementation in real time on smart phone devices, We propose a information theory based ranking of features as the preprocessing step for this purpose. In this approach the features or attributes are ranked using information gain as the criterion, and other insignificant features are discarded. This has worked surprisingly well as compared to other attribute selection methods, given that in this application context, we are dealing with very high-dimensional datasets, where we need to use around half the attributes to achieve the same level of recognition performance. We carried out extensive experiments with different features ranked by the information theory based ranking approach, including baseline traditional Naïve Bayes classifier, Decision tree, random forests, the classifiers based on ensemble learning (random committee), and lazy learning (IBk). Brief details of some of the classifiers examined for this work is given below:

#### 4.1. Naïve Bayes Classifier

This classifier is based around Bayes' theorem and computes probabilities in order to perform Bayesian inference. The simplest Bayesian method, Naive Bayes, is described as a special case of algorithm that needs no adaptation to data streams. This is because it is based on supervised learning, and it is straightforward to train the model, and performs well in terms of accuracy and generalization, making it a good method for baseline comparison. Further details of this classifier approach are given in <sup>11, 12</sup>.

#### 4.2. K-means Clustering

Clustering is an unsupervised learning approach, and here the dataset does not need to have labelled data. The instances are grouped and if they are either the same or related to each other they are placed in one group and those which are different or un-related are placed in another group. K-Means is known to be the simplest and the most popular algorithm and based on some criterion (Euclidean distance or Manhattan distance) it analyses if the instances can be clustered without having any previous knowledge about them. Due to simplicity, and its capability to work on unlabeled data, it is a good candidate for baseline reference for examining classifier performance. Further details of this classifier approach are available from <sup>13</sup>.

#### 4.3. Decision Trees

Decision Tree classifiers are based on predictive machine-learning models, that determine the dependent variable, or the target value of a new sample from the various attributes of the data available. Here, different attributes are denoted by the internal nodes of the decision tree, and the possible values that the attributes can have in the observed samples is denoted by the branches between the nodes. Further, the final values (classification) of the dependent variable are represented by the terminal nodes. The dependent variable denotes the attribute that needs to be predicted, and its value is determined by values of all other attributes. The independent variables in the dataset then form the independent attributes, and they help in predicting the value of the dependent variable. The following simple algorithmic approach is followed by the J48 Decision tree classifier used in our experiments. It has to first create a new decision tree based on the attribute values of the training data available, so as to classify a new item. The subsequent new set of items in the training set are recognized by the attributes that discriminates various samples in a clear manner. This feature provides us with most useful information, as it provides clear distinction between data instances needed for best classification. The details of J48 decision tree classifier is provided in <sup>14, 15</sup>.

#### 4.4. Random Forests

Random Forests are an ensemble of decision trees, and are based on ensemble learning methods for classification and regression. They are also thought of as form of a nearest neighbor predictor, that construct a number of decision trees at training time and output the mode of the classes as the output class. (Random Forests is invented by Leo Breiman<sup>15</sup>, and stands for an ensemble of decision trees). Random Forests try to reduce the issues with high bias and variance by computing an average, and balancing the two extremes. Moreover, Random Forests have very few parameters to tune and most of the time work very well by simply using them with parameter settings set to default values. Due to these advantages, it is often possible to use Random Forests off the shelf, without much handcrafting or modelling needed with other classifiers, and yield a reasonable model that is fast and efficient.

#### 4.5. Random Committee

Random committee is also a form of ensemble learning approach and based on the assumption of improving performance by combining classifiers. In this type of classifier, a different random number seed is used for each classifier construction; however, they are based on the same data. It then computes an average of predictions generated by each of these individual base classifiers, and outputs this average as the output class<sup>15</sup>.

#### 4.6. Lazy IBk Classifier

Lazy learners classifiers are based on the principle of learning on fly during classification time, and in fact store the training instances during training time. IBk classifier is very similar to k-nearest neighbor classifier. As most of the learning happens during classification phase, they tend to be slow, and it is possible to speed up the job of finding the nearest neighbors, by using a variety of different search algorithms. A linear search technique was used for this work, but the performance can also be enhanced by using kD-trees, or cover trees. The distance function

used was Euclidean distance. The number of neighbours used were 1, with no weighting based on distance from the test instance.

## 5. Experimental Results

To evaluate the performance of the proposed data mining approach for automatic human activity recognition from smartphone data, we used the part of the dataset <sup>1</sup>, with pre-processed feature set with 561 features and represents 2.56 seconds of time. Each record consists of following attributes:

- Triaxial acceleration from the accelerometer (total acceleration) and the estimated body acceleration.
- Triaxial Angular velocity from the gyroscope.
- A 561-feature vector with time and frequency domain variables.
- Its activity label.
- An identifier of the subject who carried out the experiment.

Num							Num						
Feat	KM	NB	J48	RF	RC	IBK	Feat	KM	NB	J48	RF	RC	IBK
2	38.00	49.45	56.30	55.60	60.10	53.18	2	15.1	0.0	0.9	7.3	14.4	0.0
8	68.40	48.26	61.39	63.01	63.03	60.18	8	20.6	0.0	7.4	16.8	17.7	0.0
16	69.00	48.57	69.02	71.27	71.10	67.84	16	37.4	0.3	11.4	19.7	23.4	0.0
32	70.00	52.34	70.24	74.17	75.10	71.74	32	67.9	0.9	25.7	25.7	25.4	0.0
64	59.00	56.10	77.30	77.51	83.73	77.51	64	119.4	1.7	38.0	29.1	31.5	0.0
128	59.50	55.31	91.46	94.29	95.10	92.97	128	217.0	4.4	64.6	31.7	30.7	0.0
256	57.00	53.86	93.81	95.63	96.28	97.55	256	457.5	3.3	52.7	20.1	25.7	0.1
561	60.00	79.00	94.00	96.30	96.90	97.89	561	582.1	5.8	247.4	14.7	27.0	0.5

(a)

(b)

Fig. 2. Comparison of Classifier Performance: (a) Recognition Accuracy; (b) Model Building Time

Figure 2 depicts the performance of different features ranked based on information gain with different classifier learning approaches, and time taken to build the model. As can be seen in this figure, we examined 2,8,16, 32, 64, 128, 256 and 561(all features), ranked by information gain approach. The total data size for training and testing comprised around 10,000 samples. We used 5 fold cross validation for partitioning this large dataset (10,000 samples) into training and testing subsets. Figure 2 shows the comparative performance for each classifier, in terms of Classification Accuracy and Time taken to build the model, and Table 1 to Table 4 shows details of other measures of performance in terms of TPR(True Positive rate), FPR (False Positive Rate), PR(Precision), RC(Recall), F-m(F-measure) and ROC. Also, the confusion matrix for best performing IBk classifier for 128 and 256 ranked features is shown in Table II.

Further, as can be seen in Figure 2, the Naïve Bayes Classifier performs reasonably well for such a large dataset, with 79% accuracy, and it is fastest in terms of building the model taking only 5.76 seconds. However, random forests, one of the ensemble learning approach is better in terms of both accuracy and model building time, with 96.3% accuracy and 14.65 seconds model building time. The other ensemble learning classifiers (random committee and random subspace), though perform well in terms of classification accuracy (~ 96%). As expected, the k-Means clustering being an unsupervised approach performs poorly with 60% classification accuracy, and 582 seconds. The best performing classifier however, is IBk classifier, which is based on lazy learning, resulting in an accuracy of more than 90% for 128 features and 256 features.. A trade-off between accuracy and model building time is necessary for a smartphone based activity recognition system, as real time activity monitoring needs to model to be built dynamically from the captured data, and faster model building time with accuracy recognition accuracy is the best to aim for. Further, additional performance measures such as TPR, FPR, Precision, Recall, F-measure and

ROC area need to be taken into consideration for choice of best algorithm for building automatic activity recognition systems. Table I shows these additional performance measures.

Table 1. Additional Performance Measures For Different Classifier Models

Features	TPR	FPR	PR	RC	F-m	ROC
IBK(256)	0.976	0.005	0.976	0.976	0.976	0.985
RC(256)	0.963	0.008	0.963	0.963	0.963	0.998
RF(256)	0.956	0.009	0.956	0.956	0.956	0.998
RC(128)	0.951	0.01	0.951	0.951	0.951	0.996
RF(128)	0.943	0.012	0.943	0.943	0.943	0.996
J48(256)	0.938	0.012	0.938	0.938	0.938	0.971
IBK(128)	0.93	0.015	0.93	0.93	0.93	0.957
J48(128)	0.915	0.017	0.915	0.915	0.915	0.96
RF(64)	0.837	0.035	0.839	0.837	0.838	0.969
RC(64)	0.837	0.035	0.839	0.837	0.838	0.969
IBK(64)	0.775	0.049	0.776	0.775	0.776	0.863
J48(64)	0.773	0.048	0.774	0.773	0.774	0.891
RC(32)	0.751	0.053	0.755	0.751	0.752	0.943
RF(32)	0.742	0.055	0.746	0.742	0.743	0.942
IBK(32)	0.717	0.06	0.72	0.717	0.718	0.826
RF(16)	0.713	0.061	0.716	0.713	0.714	0.929
RC(16)	0.711	0.061	0.714	0.711	0.712	0.923
J48(32)	0.702	0.063	0.701	0.702	0.701	0.857
J48(16)	0.69	0.065	0.69	0.69	0.689	0.873
IBK(16)	0.678	0.067	0.681	0.678	0.679	0.812
RF(8)	0.63	0.079	0.63	0.63	0.63	0.889
RC(8)	0.63	0.078	0.633	0.63	0.631	0.872
J48(8)	0.614	0.082	0.61	0.614	0.611	0.868
IBK(8)	0.602	0.084	0.606	0.602	0.603	0.791
J48(2)	0.563	0.092	0.565	0.563	0.562	0.888
NB(64)	0.561	0.093	0.57	0.561	0.515	0.878
RF(2)	0.556	0.093	0.558	0.556	0.557	0.876
NB(128)	0.553	0.09	0.644	0.553	0.523	0.928
NB(256)	0.539	0.092	0.662	0.539	0.505	0.928
IBK(2)	0.532	0.097	0.539	0.532	0.533	0.854
NB(32)	0.523	0.102	0.517	0.523	0.453	0.873
NB(2)	0.495	0.108	0.512	0.495	0.418	0.865
NB(16)	0.486	0.11	0.494	0.486	0.417	0.86
NB(8)	0.483	0.11	0.489	0.483	0.413	0.859

Table 2. Confusion Matrix for IBk Classifier for 128 and 256 ranked features (best performing classifier)

	RC-0	RC-1	RC-2	RC-3	RC-4	RC-5
AC-0	1717	3	2	0	0	0
AC-1	12	1513	19	0	0	0
AC-2	9	25	1372	0	0	0
AC-3	0	0	0	1471	235	71
AC-4	0	0	0	231	1658	17
AC-5	0	0	0	71	28	1845

The confusion matrix for IBk classifier (the best performing classifier) for 128 and 256 ranked features is shown in Table II. The confusion matrix shows how the classifier confuses and misclassifies one class for another (Actual



Class (AC-0 to AC-5) vs. Recognised Class (RC-0 to RC-5). As can be seen in Table 2, the classifier performs has least confusion in recognising the Class 6 (Laying) and Class 1 (walking activity) out of 6 different activities (WALKING, WALKING\_UPSTAIRS, WALKING\_DOWNSTAIRS, SITTING, STANDING, LAYING). It confuses little bit between walking upstairs and walking downstairs activities and confuses bit more between sitting and standing. With just one smartphone tied to the waist, this is a significantly better performance for recognising each activity.

## 6. Conclusions And Further Plan

In this article, we proposed a novel automatic activity recognition scheme using smartphone data based on optimal attribute selection based on information theory based ranking and machine learning techniques. We examined several learning approaches and found lazy learning, random forests and ensemble learning based approaches to be promising in terms of activity classification accuracy, model building time for automatic classification, and confusion matrix, with experimental validation on publicly available activity recognition dataset. Further research would involve adapting the proposed data mining approach on embedded hardware using appropriate implementation process, so that it can be implemented on smartphone devices. Also, other active and novel unsupervised learning approaches need to be investigated as model building in real time on resource constrained smartphones could be restrictive.

## References

1. Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra and Jorge L. Reyes-Ortiz. Human Activity Recognition on Smartphones using a Multiclass Hardware-Friendly Support Vector Machine. International Workshop of Ambient Assisted Living (IWAAL 2012). Vitoria-Gasteiz, Spain. Dec 2012.
2. Ekholm, J., Fabre, S.: Forecast: Mobile data track and revenue, worldwide, 2010- 2015. In: Gartner Mobile Communications Worldwide. July 2011.
3. Cook, D.J., Das, S.K.: Pervasive computing at scale: Transforming the state of the art. *Pervasive and Mobile Computing* 8(1) February 2012, 22-35.
4. Allen, F.R., Ambikairajah, E., Lovell, N.H., Celler, B.G.: Classification of a known sequence of motions and postures from accelerometry data using adapted Gaussian mixture models. *Physiological Measurement* 27(10)
5. Rodriguez-Moliner, A., Perez-Martinez, D., Sam\_a, A., Sanz, P., Calopa, M., Galvez, C., Perez-Lopez, C., Romagosa, J., Catalia, A.: Detection of gait parameters, bradykinesia and falls in patients with parkinson's disease by using a unique triaxial accelerometer. World Parkinson Congress, Glasgow ,2007.
6. Mannini, A., Sabatini, A.M.: Machine learning methods for classifying human physical activity from on-body accelerometers. *Sensors* 10(2) (2010) 1154{117}
7. Ravi, N., D, N., Mysore, P., Littman, M.L.: Activity recognition from accelerometer data. In: In Proceedings of the Seventeenth Conference on Innovative Applications of Artificial Intelligence (IAAI, AAAI Press, 2005, 1541-1546.
8. Kwapisz, J.R., Weiss, G.M., Moore, S.A.: Activity recognition using cell phone accelerometers. *SIGKDD Explor. Newsl.* 12(2), March 2011, 74-82.
9. LeCun, Y., Jackel, L., Bottou, L., Brunot, A., Cortes, C., Denker, J., Drucker, H., Guyon, I., Mller, U., Sckinger, E., Simard, P., Vapnik, V.: Comparison of learning algorithms for handwritten digit recognition. In: International Conference on Artificial Neural Networks. 1995, 53-60.
10. Ganapathiraju, A., Hamaker, J., Picone, J.: Applications of support vector machines to speech recognition. *Signal Processing, IEEE Transactions on* 52(8) (aug. 2004) 2348-2355.
11. SME Hossain and G. Chetty. Next Generation Identity Verification Based on Face-Gait Biometrics. In The International Conference on Biomedical Engineering and Technology, ICBET '11, 2011.
12. G. Chetty and M. Wagner, Investigating feature-level fusion for checking liveness in face-voice authentication, Eighth International Symposium on Signal Processing and Its Applications, August 28-31, 2005, p. 66-69..
13. Ashraf, Mohammad, GirijaChetty, Dat Tran, and Dharmendra Sharma. "A New Approach for Constructing Missing Features Values." *International Journal of Intelligent Information Processing* 3, no. 1, 2012.
14. Leo Breiman, Random Forests, *Springer Machine Learning Journal*, Oct.. 2001, Volume 45, Issue 1, p. 5-32.