# User Exercise Pattern Prediction through Mobile Sensing

Georgi Kotsev, Le T. Nguyen, Ming Zeng, and Joy Zhang
Carnegie Mellon University Silicon Valley
Moffet Field, California, USA
{georgi.kotsev, le.nguyen, ming.zeng, joy.zhang}@sv.cmu.edu

*Abstract*—**Even though the health benefits of regular exercising are well known, an average person has difficulty maintaining physical activity on a regular basis. One of the main reasons for this is lack of motivation. With their increasing ubiquity, wireless devices and smartphones and their sensing capabilities now can be involved in solving this issue. Many mobile applications have been developed with which people are able to keep track of their exercises, become more aware of their physical condition, and be more motivated. The collected data is also a good source for researchers in understanding the exercise patterns and the main factors influencing people to exercise. Understanding those factors will allow better applications to be built, which helps motivate people. In this work, we quantitatively analyze a dataset collected from over 10,000 users. To better understand the user exercise patterns, we identify a set of factors influencing their exercise patterns. Based on these insights, we develop a prediction model to predict users' future exercise activities.**

## I. INTRODUCTION

Regular physical activity does not only improve people's physical condition but also their mental health [9]. However, the majority of adult Americans are not physically active on a regular basis, which increases individual's risk of diseases [9].



Fig. 1.   Wireless devices for fitness tracking

Even though people are aware of the fact that exercising can improve their health, there are numbers of factors that influence both their motivation and their exercise behavior. In order to motivate them successfully, we need to identify and understand these factors first, and then we will be able to use that information to intervene in peoples' exercise behavior if needed.

One intuitive factor influencing the exercise habits could be the weather. If it is cold outside or raining, some people prefer to stay at home and postpone their workout for another day. In Figure 2 we can see the average activities per month for users in New York City and the development of the temperature
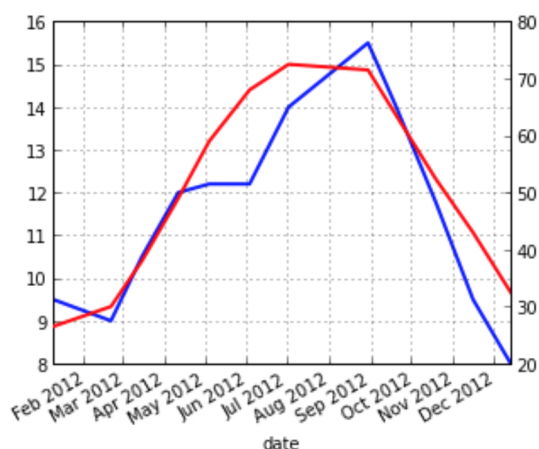


Fig. 2.   Average exercises per month in blue and average temerature in red for New York City in 2012. The left Y-axis shows the average exercises and the right Y-axis shows the temperature in Fahrenheit

over the same year. By observing both we can find a clear correlation between the two curves. Compared with other U.S. major cities, New York City has distinguishable seasons, and by looking at the curves, it seems that the seasons correlate with the activity level of the users. However, there is an intersting fact - although the temperature is colder in January than in December - we observe that starting from January the activity level increases. This may be affected by the New Year's resolution of people commiting to a healthier lifestyle. This suggests that commitment is another factor that makes people exercise more.

Humans are social by nature. This social factor could also be of high significance in determining why people exercise. Some people find it much easier to run with a running buddy or in a small group. When they are part of a team, they like to keep up with the other members, and they exercise even on days that they are not motivated.

An interesting research question would be what are the other factors? Especially, what are the other indicators that would help us predict users' future exercising habits? Also, what is the level of commitment of a user, and can you tell the length of the commitment by just looking at the beginning period? Commitment in this sense means the number of exercises for a period in the future.

In recent years the wireless fitness tracking devices are getting more and more popular. Gadgets like Fitbit, Nike+, and Jawbone are becoming a necessity for the modern passionate jogger. On the other hand, many fitness tracking applications for smartphones have been developed lately. They use the sensing capabilities of the smartphones to provide similar functionality as the dedicated devices. Using these devices not only keeps the users motivated to exercise but also helps them understand their exercise patterns. Applications, such as Endomondo, Runkeeper, or Strava, leverage the use of mobile phones to capture users' exercise metrics, including speed, total amount of miles, and an estimated amount of burned calories.

For the first time in the history we are able to collect large amounts of data from the tracking applications. Without mobile devices, collecting all this data would be extremely difficult, and the scaling would be with a huge overhead. By looking deeply into the data and discovering patterns, we can understand human exercise behavior. Analyzing this data helps us infer the factors influencing people to exercise. This will allow us to motivate users to exercise more when we detect a need.

In this work, we present insights about user exercise patterns by analyzing historical data of over 10,000 participants. Our goal is to understand the factors that influence the user's activity patterns. We create a model to predict whether the number of exercises increases or decreases in the future. Our contribution in this work can be summarized as follows:

- **Insights about users' exercise patterns**: We introduce insights about the fitness behavior of the users.

- **Influential factors**: We analyze the important factors for the understanding of users' exercise patterns.

- **Prediction**: We propose a modeling approach to predict the tendency of users' future number of exercises per week and compare the performance of different predictors and classifiers.

By understanding the factors that impact user's motivation and behavior, we can build systems that better adapt to a user's individual profile and needs. For example, if a decrease of exercise frequency can be predicted, an assistant system can help motivate the users through reminders to keep exercising. Compared to taking medications for increasing health, maintaining a healthy lifestyle and performing regular fitness activities are better alternatives.

In the next sections of this paper we describe the data set and provide surface statistics about the user in order to understand them better. Then we present different variations of a prediction model and compare their performance with different parameters. We conclude with a presentation of the related work and discuss possible future works.

## II. DATA DESCRIPTION

### A. Nature of the Data

There are many mobile applications used for tracking users' exercises. For example, for tracking running exercises, a user can use an mobile application, which leverages the GPS to track his location, Speed can also read data from different

| User profile information | Exercising information (collected by mobile sensors) |
|---|---|
| total number of exercises | activity id |
| total distance done | user id |
| total calories burned | date |
| location | distance |
| gender | duration |
| motivation | activity type |
| months active | location |
| number of friends | |

TABLE I.    TYPES OF DATA

additional sensors, such as a heart rate monitor, if available. During exercise the application would provide monitoring information about the distance, time, and speed. After the workout the user usually uploads his workout in the cloud, where he can keep track of his own activities, share them with other users, and also see others' activities and achievements.

Users would also have the opportunity to post their workouts on social media such as Facebook or Twitter. Typically there are two types of data, one being profile data, static with basic information about the user. The other being time series data, which describes every exercise that the user did with date and duration - both can be seen in Table I. With the use of such an application, we were able to collect data for over 10,000 users, who are described in the next subsection.

### B. User Demographics

*1) Activity Type:* The user has the opportunity to track different activities. Besides running, users can keep track of a variety of activities, including swimming, housework, or any user-defined activities. The percentage distribution of the top 10 is shown in Figure 3. The top 3 most popular activity types are running, walking, and cycling.
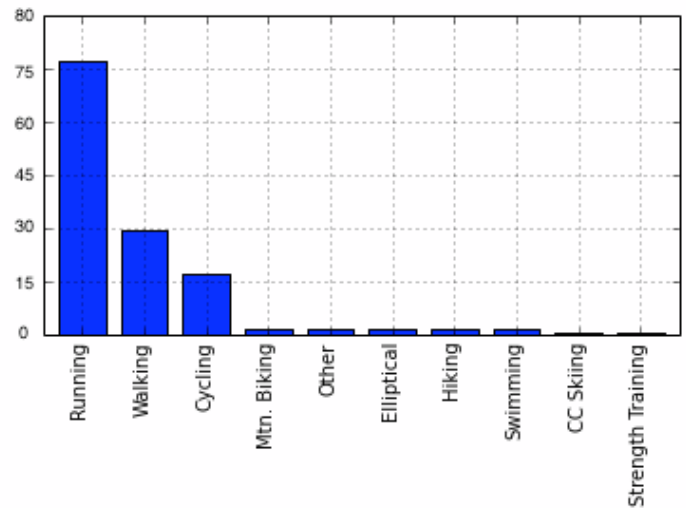


Fig. 3.    Percentage of the user doing each activity type for the top 10 most popular activities. "Other" is a label of activity entered by the user. A user can keep track of multiple activities

From the given data, it can be concluded that a person uses the application to keep track of solely one or multiple activity types. By computing the activity entropy for each user, we can find out if one uses the applicatoin to keep track of single or multiple activities.

*2) Location:* The majority of the users are in the US, followed by Sweden and the Netherlands. This can be seen in Figure 4. One interesting fact is that even though the UK has more than six times the population, we have more than two times more Swedish users than British. However, we could not draw the conclusion, that the Swedish are 12 times more sporty than British people as there could be number of reasons for that.
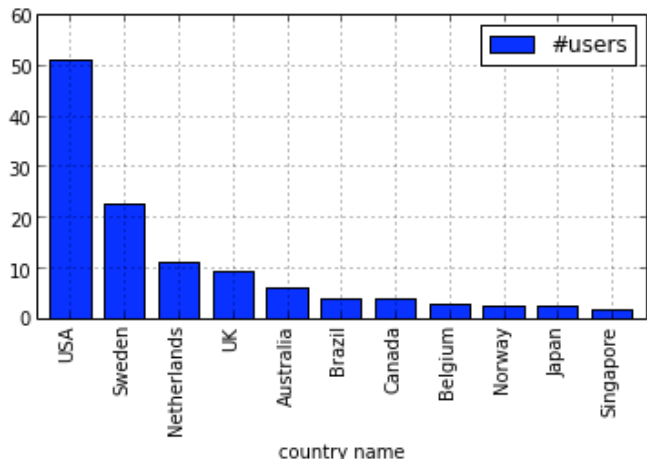


Fig. 4.   Percentage of users in top 10 countries.

Studying the US users, we can see that the state with the most users is California, followed by Texas and Florida. California is also the most populated state, so this is more or less expected. A more interesting fact is that the State of New York and Texas have similar populations, but we have two times more users from Texas than in New York state.
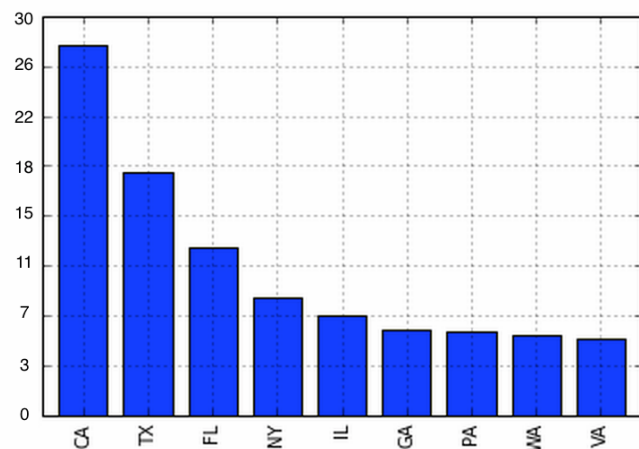


Fig. 5.   Percentage distribution of all users in top 10 states.

If we take a closer look at the runners in the state of California, we can see that although San Francisco is not the biggest city, it has more users who try to be aware of their health and keep track of their fitness exercises, followed by cities like Los Angeles and San Diego.

*3) Gender:* The gender distribution of the users and their exercise frequency can be seen at Table II. Here the data shows
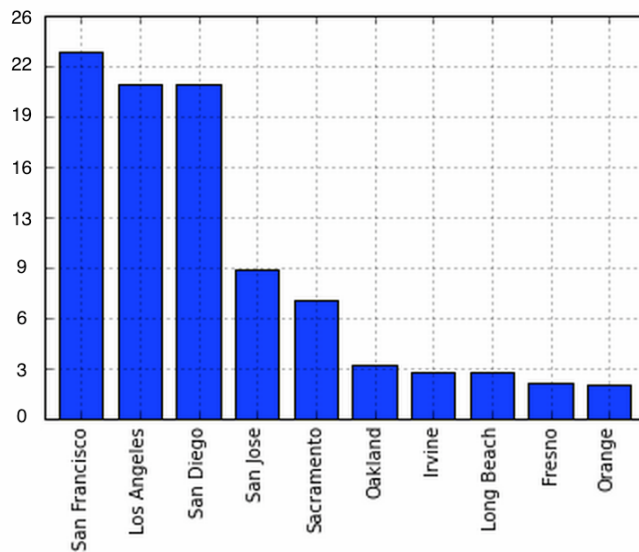


Fig. 6.   Percentage distribution of the Californian users in top 10 cities.

| Activities / Month | Male | Female |
|---|---|---|
| User Percentage | 58 | 42 |
| Mean | 7.06 | 7.86 |
| Standard Divination | 4.48 | 4.70 |

TABLE II.   GENDER AND EXERCISE FREQUENCY

that we have more male users, but the female users tend to exercise a little bit more often than males.

As we can see from Table III, women walk longer distances than men, but men run and cycle longer distances than women.

As we can also see from Table IV, women spend more hours walking, but men tend to spend more time running and cycling.

*4) Social Connections:* In the application you have also the opportunity to connect with other users by adding them as friends. In Figure 7 we can see the distribution of the friend-group sizes. In this distribution we can observe that the majority of the users have less than 10 friends.

Out of this research comes the interesting question is there a correlation between the numbers of friends of a users and how active he/she is? In Figure 7 we can observe the percentage distribution of the users. Definitively the majority of the the users have less then 10 friends.

| Activities / Month | Male - Mean (SD) | Female - Mean (SD) |
|---|---|---|
| km walked / Month | 5.12 (10.63) | 6.94 (11.88) |
| km run / Month | 14.37 (17.05) | 10.88 (14.07) |
| km cycled / Month | 7.38 (10.6) | 4.73 (8.11) |

TABLE III.   GENDER AND DISTANCE

| Activities / Month | Male - Mean (SD) | Female - Mean (SD) |
|---|---|---|
| Hours / Month | 3.59 (3.49) | 3.79 (3.81) |
| Hours walked / Month | 1.42 (2.76) | 1.95 (3.06) |
| Hours run / Month | 2.32 (2.64) | 2.22 (2.75) |
| Hours cycled / Month | 0.68 (1.13) | 0.49 (0.93) |

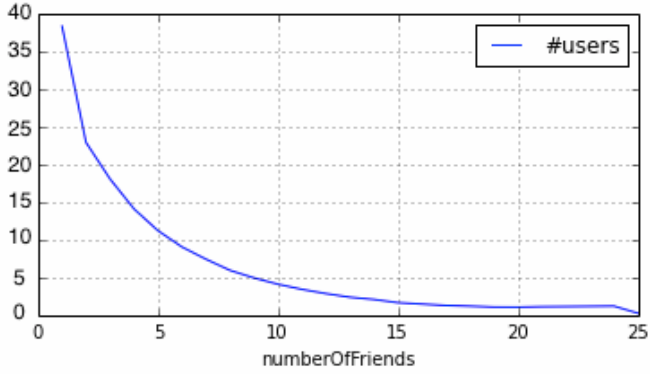TABLE IV.   GENDER AND TIME DEDICATION

Fig. 7. Percatage of users having a certain number of friends. Majority of the user have under 10 friends.

*5) Time: Day of the week*

In Figure 8 we can see how the distribution of the weekly runs looks like. For each user we have calculated what percentage of his/her workouts are on the given day. In the figure we present a box plot of that data. From the graphic we can say that Tuesday leads the ranking of favorite days of the week for exercising. On the weekends the user may have activities other than sports, and on Monday they start their work week. So Tuesday begins most people's weekly workout schedule.
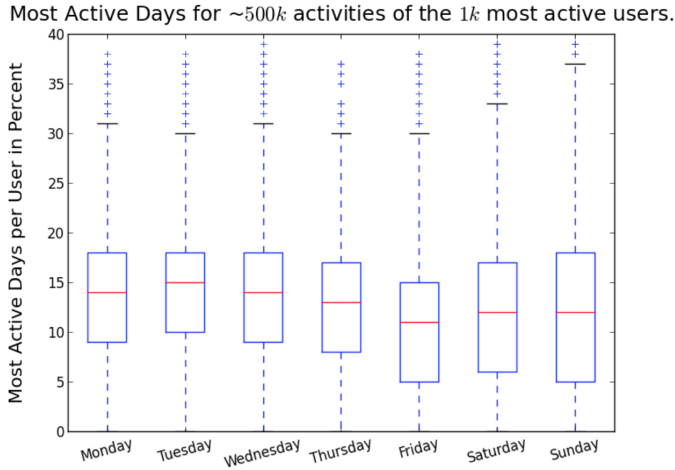


Fig. 8. Days of the week distribution.

*Month of the year*

The month of the year could be also form significance for the number of activities. Figure 9 shows the average of activities per month for the top 10% most active users. We can observe the intuitive seasonal influence of the exercise habits of the users. In the warm summer months they tend to go out and run more compared to the colder winter months.

In Figure 10 we can see the average exercise frequency of all users, clustered by countries. We can see that American, Canadian, and Japanese users exercise the more frequently than the rest.
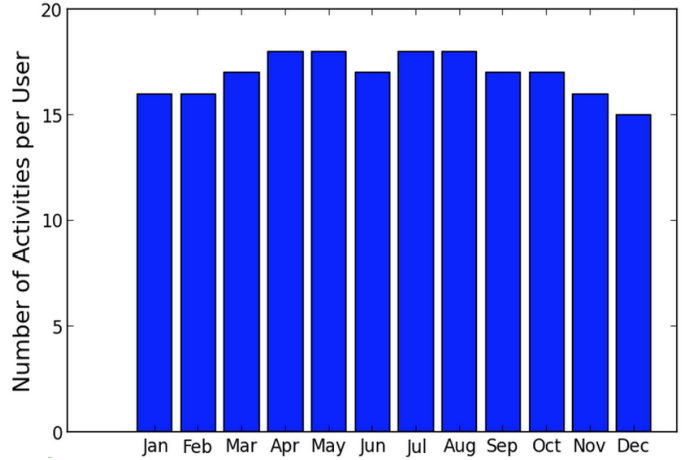


Fig. 9. Distribution of average number of activities per months for the top 10% most active users
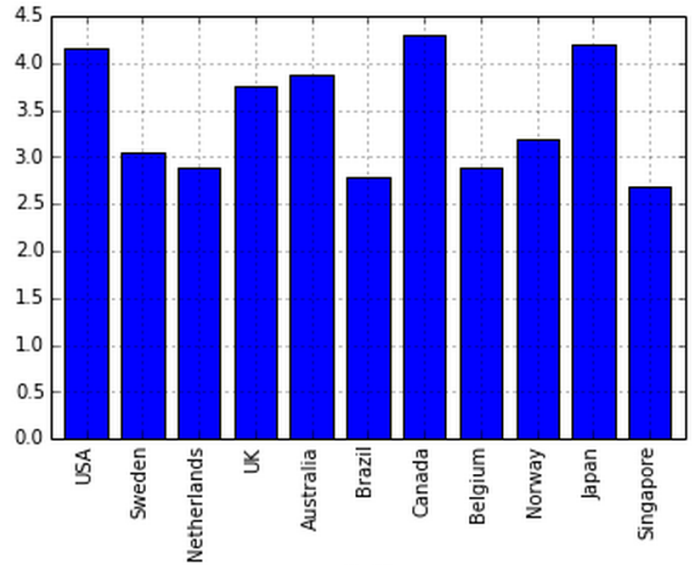


Fig. 10. Average exercise frequency per month in different locations for all users

## III. PREDICTION MODEL

In this section we present the quantitative analysis of some of the main factors influencing people's exercise patterns and compare the performance of different prediction models. For the purpose of the prediction, we focus our further analysis only on the running activities.

*A. Problem Statement*

In order to predict the future activities of our users, we test the performance of several prediction methods. The first one is mathematical, only using the tendency of recent time. The other employs different machine learning estimators to learn the features and make a prediction.

In Figure 11 we observe a sample time series data from one user. It presents the number of exercises that a person performs a week. We can see that at his/her first months

he/she was active at some point, but then in the next period he became not that dedicated to the sport. It would be interesting to discover what made him behave this way. Was it just the seasonal change from summer to winter? What if there are some specific trait of his/her profile that makes it possible to predict these patterns?
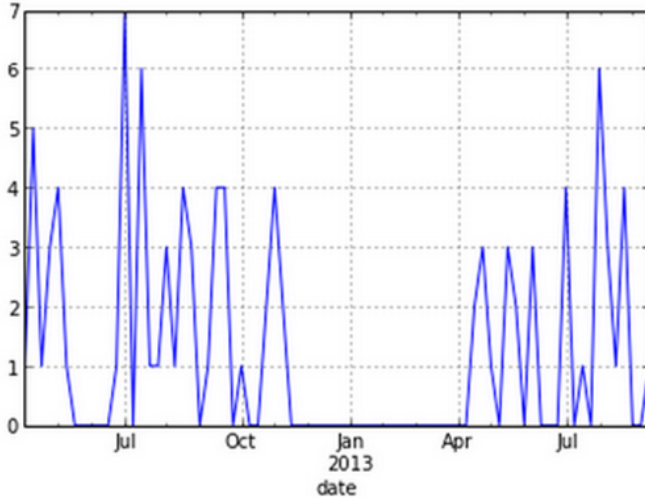


Fig. 11. Runs per week time series for user A.

We can analyse time series with different metrics. In Figure 12 we also see an example of a time seris based on the distance. In Figure 13 we can see the speed of the same user. What is interesting is that in the beginning of 2012, his/her distance was increasing, but the speed decreased respectively.
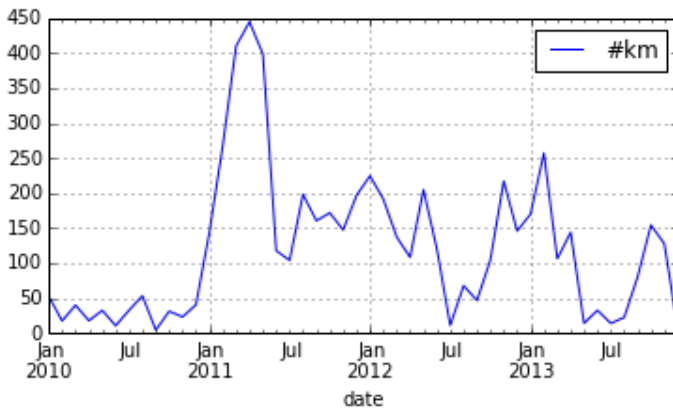


Fig. 12. Distance time series of user B

In this research we are trying to answer the question of if and how we can combine this time series to predict the future values of the number of exercises.

First we look at a tendency-based approach, which is basicly one feature time series analysis. Then we propose a machine learning model, which is based on the initial period of time that can make predictions for the future performance of the user. With learning the generated features by different users in their first period we train the model to predict the change of the features, e.g. comparing the average speed of the runs in the first period, and the average speed in the second period.
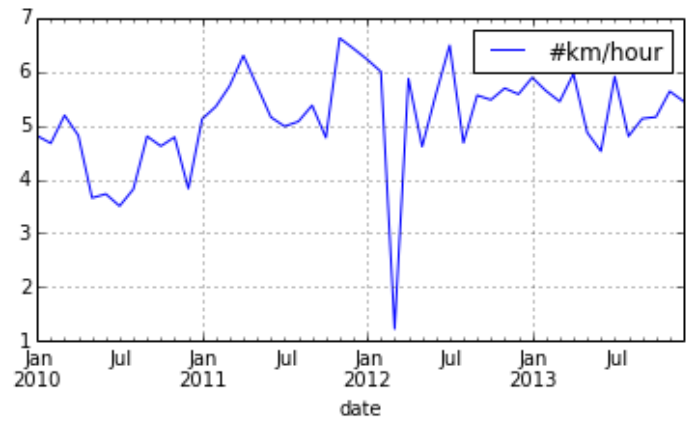


Fig. 13. Speed time series of user B

Looking at the initial period, we define different classes describing the change of the number of exercises in the next period. We enumerate three classes with -1, 0 and 1 where 1 means increase, -1 means decrease, and 0 means there is no significant change. We also choose the threshold for a significant change to be around 30%, because this threshold divides the classes most evenly.

### B. Tendency-Based Prediction

We could try to predict all of the three time series - runs, distance, or speed. Of the three, the most helpful would be the frequency of runs a week, because if we detect that the tendency is decreasing, we we can trigger some application functionality to motivate people to continue running.

$$P(a_{t+1}) = tendency(a_{t+1} - a_t) = tendency(a_t - a_{t-1}) \quad (1)$$

In Equation 1 we could see the proposed approach. The predicted future tendency for time moment *(t+1)* is the current one. Tendency is the change of *a*, the number of activities. It could be *increasing*, *decreasing* or *level*. This approach results in average of 70% accuracy.

### C. Machine Learning Prediction

For the Machine Learning approach we frame the problem as a classification problem. We have experimented with different targets, and the results showed that to predict the number of exercises per week of the user is a reasonable metric. This way we can detect if a user is going to exercise less and a scenario of a motivational application that gives him/her an incentive to keep on exercising.

In Table V we can see one sample training vector. For example, in this case we have a user running two times in the first week with mean runs per week for the initial period of 2.5 and also increasing his/her runs per week in the future, according to the target label. The target label represents the class to which his behavior is classified and gives us information if he/she is going to decrease, increase, or stay with the same number of exercises per week.

| feature | value |
|---|---|
| runs first week | 2 |
| mean runs per week | 2.5 |
| max runs per week | 3 |
| min runs per week | 2 |
| ave speed first week | 3.86 |
| .. | . |
| .. | . |
| target class | 1 |

TABLE V.     SAMPLE TRAINING DATA ENTRY

| Index | Classifier |
|---|---|
| 1 | Base Line |
| 2 | Gaussian Naive Bayes |
| 3 | Decision Tree |
| 4 | Random Forest |
| 5 | Extra Trees |
| 6 | Nearest Centroid |
| 7 | K-Neighbors |
| 8 | Logistic Regression |
| 9 | Linear SVC |
| 10 | SVC |
| 11 | Linear Regression |

TABLE VII.     CLASSIFIERS USED

*1) Features Extraction:* Apart from the time series data that we have for each user for the different measures, number of runs, distance, duration, and speed. We extract also additional statistical features for each user such as mean, standard deviation, variance, min, max and range. These data sets will be helpful, because they summarize the performance of the user for the period.

Having more features does not necessarily mean achieving higher accuracy. Some of them could contain redundant or irrelevant information. This can cause adding more noise to the classifier and result in a decrease of accuracy. The informative features are those, which make the data easily separable that will result in being more predictable.

For determining the best features, we use a univariate feature selection with the F-test. Each feature receives a score of relevancy. The higher the score, the more relevant the feature is. Table VI shows top five features form all 50 features that we have considered.

| Feature | P-Value | Score |
|---|---|---|
| runs first week | 6.8e-12 | 26.29 |
| mean runs per week | 2.1e-08 | 17.94 |
| max runs per week | 4.4e-07 | 14.81 |
| min runs per week | 2.2e-06 | 13.16 |
| ave speed first week | 1.4e-04 | 8.93 |
| 2 elapsed hours | 1.1e-03 | 6.81 |
| min distance | 1.2e-03 | 6.73 |
| min elapsed hours | 4.8e-03 | 5.37 |
| mean speed | 2.6e-02 | 3.68 |
| max speed | 2.6e-02 | 3.66 |

TABLE VI.     TOP 10 FEATURE RANKING, BASED ON THE SINGLE USER MODEL

*2) Single User Model:* For the single user model we train a separate classifier for each user. The classifier learns by the historical data of the user. This is basically a time series prediction problem.

We aggregated the number of activities that each user has for one week and presented the data as time series. In Figure 11 we can see the exercise patterns of one user. In order to predict the change of the number of activities, we could use the features that we have from the historical exercises data. With the information we have for the particular users, we could also infer additional statistical features that could help our prediction be better.

In Table VII we can see the classifiers that we have employed for the prediction. We have also performed a 10-fold cross validation. The results can be seen in Figure 14. As we see from the results, with Support Vector Machines we are able to achieve accuracy of 72%.
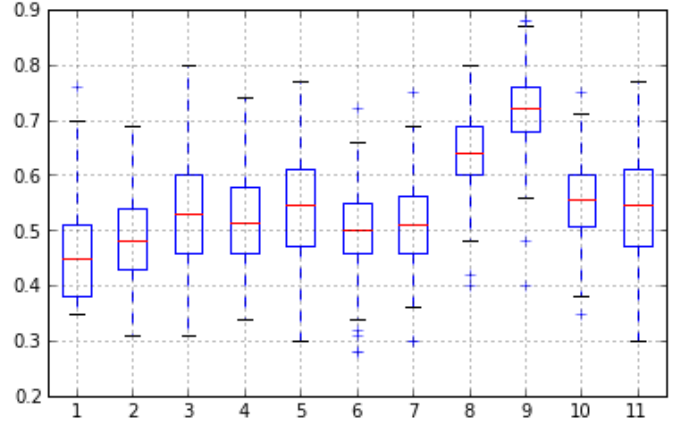


Fig. 14.   Single user model classifier results

*3) Multi User Model:* In contrast to the single user model where for each user we train different predictive model, for the multi user model, we train the same model with all of the users. This way we can use the collective history in order to predict the future behavior for a single user. More specifically, we look at an initial period (e.g. first two weeks of usage of the application) and then try to predict if the future performance, number of activities, or exercising frequency are going up, down, or stay level.

We used a data set of 10,000 users to train 10 different classifiers and trained them with the generated features. We performed a 10-fold cross validation for each classifier. Table VIII shows the mean accuracy for the top three best performing classifiers and the base line case. The base line is a naive predictor for classification using the most frequent class.

The quality of the prediction depends on the length of the initial period that we are analyzing and extracting our features from. We also quantitatively analyze the influence of changing this window size of the initial period and present the results in Figure 15. From our data we infer that looking in the first three weeks of users' exercise patterns is relatively good enough for predicting the future number of runs of a user.

## IV.   RELATED WORK

There is a variety of mobile device applications that are taking advantage of their sensor capabilities to collect detailed

| Classifier | Accuracy rate in % |
|---|---|
| Base Line | 0.39 |
| Decision Tree | 0.45 |
| Extra Trees | 0.47 |
| Random Forest | 0.49 |
| K-Neighbors | 0.51 |
| Logistic Regression | 0.53 |

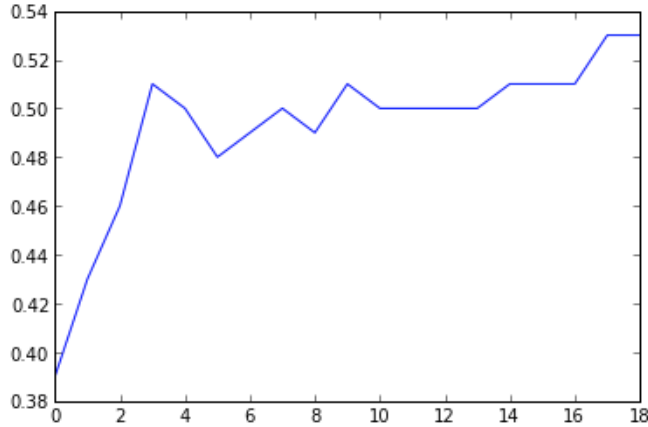TABLE VIII.    MULTI USER TOP FIVE CLASSIFIER RESULTS



Fig. 15.    Change in accuracy depanding on the window size

data, which can be used for building statistical models for a better understanding and prediction. One example is the Ladha et al. publication [7] from 2013. The authors leverage the use of mobile sensors to collect accelerometer data from the climbing domain. They present an approach for climbing performance analysis as a replacement of expert assessments. They evaluated their results with climbers in a competition and suggest that their work represents a first step towards an automatic coaching system for climbing enthusiasts.

In another publicaiton, the authors try to motivate physical activity at work, using persuasive social media and mobile devices [3]. Tring to increase physical activity, using pedometers is also a topic covered in [1] and [4].

The authors of [8] analyze the trends of mobile-health applications for mobile devices. They classify the top 200 apps and show that although the biggest group of apps were delivered from or related to medical articles, websites, or journals, mobile users disproportionally favored tracking tools. In [2], [12], and [11] the authors also analyze how the physical activity increases with employment of mobile devices and social media. In [10], [6], and [5], the authors discuss the usability mobile applictions for weight loss as persional fitness coach. They collected qualitative feedback in a user study, which showed good evaluation results.

## V.    CONCLUSION AND FUTURE WORK

In this paper we proposed an innovative approach, contributing to understanding and solving a problem that is important for a significant portion of the society: lack of motivation for physical activity. We analyzed the factors that make people who exercise more different from those who exercise less. Our data set was collected by ubiquitous devices such as smartphones, whose users have shared their information online.

We framed a prediction problem as classification and tested the performance with 10 different classifiers for predicting the average future performance of a user based on information from his inital period. By generating the different features and giving them a score, we identified the most significant factors, which infer the future improvement.

In our future work, more complex models could be employed, which incorporate different information and aspects. Also the trend-based approach could be combined with the machine learing based on the confidence score. Additionally we could create a clusters or profiles of different users based on static features of them and train different classifiers for each cluster. This way the accuracy of the prediction of the model could rise. All these new strategies could potencially increase the accuracy and allow the model to do better predictions.

## REFERENCES

[1]  D. M. Bravata, C. Smith-Spangler, V. Sundaram, A. L. Gienger, N. Lin, R. Lewis, C. D. Stave, I. Olkin, and J. R. Sirard. Using pedometers to increase physical activity and improve health: a systematic review. *JAMA*, 298(1):2296–2304, 2007.

[2]  J. Fanning, P. S. Mullen, and E. McAuley. Increasing physical activity with mobile devices: A meta-analysis. *J Med Internet Res*, 14(6):e161, Nov 2012.

[3]  D. Foster, C. Linehan, S. Lawson, et al. Motivating physical activity at work: using persuasive social media extensions for simple mobile devices. 2010.

[4]  L. A. Kaminsky, J. Jones, K. Riggin, and S. J. Strath. A pedometer-based physical activity intervention for patients entering a maintenance cardiac rehabilitation program: a pilot study. *Cardiovascular Diagnosis and Therapy*, 3(2), 2013.

[5]  P. Klasnja and W. Pratt. Methodological review: Healthcare in the pocket: Mapping the space of mobile-phone health interventions. *J. of Biomedical Informatics*, 45(1):184–198, Feb. 2012.

[6]  M. Kranz, A. MöLler, N. Hammerla, S. Diewald, T. PlöTz, P. Olivier, and L. Roalter. The mobile fitness coach: Towards individualized skill assessment using personalized mobile devices. *Pervasive Mob. Comput.*, 9(2):203–215, Apr. 2013.

[7]  C. Ladha, N. Y. Hammerla, P. Olivier, and T. Pltz. Climbax: Skill assessment for climbing enthusiasts. 2013.

[8]  C. Liu, Q. Zhu, K. A. Holroyd, and E. K. Seng. Status and trends of mobile-health applications for ios devices: A developer's perspective. *J. Syst. Softw.*, 84(11):2022–2033, Nov. 2011.

[9]  A. F. Manley. Physical activity and health. a report of the surgeon general, executive summary. Technical report, US Public Health Service, 1999.

[10]  S. L. Mansar, S. Jariwala, M. Shahzad, A. Anggraini, N. Behih, and A. AlZeyara. A usability testing experiment for a localized weight loss mobile application. *Procedia Technology*, 5(0):839 – 848, 2012. 4th Conference of {ENTERprise} Information Systems aligning technology, organizations and people (CENTERIS 2012).

[11]  J. Stragier and P. Mechant. Mobile fitness apps for promoting physical activity on twitter: the #runkeeper case. In *Etmaal van de communicatiewetenschappen, Proceedings*, page 8, 2013.

[12]  D. Thompson, D. Cantu, R. Bhatt, T. Baranowski, W. Rodgers, R. Jago, B. Anderson, Y. Liu, A. J. Mendoza, R. Tapia, and R. Buday. Texting to increase physical activity among teenagers (txt me!): Rationale, design, and methods proposal. *JMIR Res Protoc*, 3(1):e14, Mar 2014.