

# Altrumetrics: Inferring Altruism Propensity Based on Mobile Phone Use Patterns

Ghassan F. Bati, *Student Member, IEEE* and Vivek K. Singh, *Member, IEEE*

**Abstract**— Altruism, i.e. *an act that one does at their own expense that tends to enhance others well-being*, is a fundamental human behavior with implications for personal and societal welfare. Hence, modeling altruism is an important building block in designing human-centered computing systems. Traditional methods for understanding an individual's altruistic propensities have been surveys and lab experiments. However, the emerging “personal big data” coming from mobile and ubiquitous devices allows for creation of lower-cost, quicker, automated methods for modeling human behaviors and propensities. We propose a new methodology to model altruism using phone data. Based on analysis of data from a 10-week field study (N=55 participants), we report that: (1) multiple phone-based features are associated with users' altruistic propensities; (2) phone features-based altruism prediction model yielded significantly better performance than a demography-based model. The results pave way for utilizing “personal big data” to model altruism in multiple commercial and social applications.

**Index Terms**—C.3.h Ubiquitous computing; O.8.15 Social science methods or tools

## 1 INTRODUCTION

With the growth in personal, mobile, and ubiquitous computing, increasingly larger aspects of human life are mediated by devices. Consequently, the data captured by such devices, i.e. the “personal big data”, is enabling a datafication of the human life and facilitating the creation of rich composite personas of different users [1]. With 1.4 billion smartphones and millions of quantified-self device users, more and more users keep track of their behavior, which has already shown value in the fields of healthcare, well-being, and urban planning [2], [3], [4], [5].

At the same time, with the growth in mobile social networks, social internet of things, and cyber-physical-social systems, there is an ever growing need to model and understand human beings as they interact with other humans and socio-technical ecosystems. For example, humans in such systems may choose to act selfishly or behave **altruistically**. For example, they may choose to provide bandwidth and resources, contribute open-source code, and write Wikipedia entries. In the emerging socio-technical landscape, they may also choose to act differently in the shared economy settings (e.g. Task-rabbit, Uber), and have different preferences for their autonomous cars and bots to behave with others. Further, altruism has been connected with emotional, physical, and financial well-being of individuals as well as communities [6], [7]. Such altruism preferences are hence at the core of the design of human-centered computing systems and their importance is only going to increase in the coming years.

In this paper, *we define altruism as any act (or behavior) one does at their own expense that tends to enhance others well-being*

[6]. There have been multiple efforts aimed at trying to elicit an individual's propensity to behave altruistically with others [6], [8]. However, previous works have generally focused on traits that could be simply observed (e.g., gender, ethnicity, age) or elicited in a small period in lab settings (e.g., via surveys and games). Regrettably, the human-related information taken by observations in restricted and atypical settings involved a limited amount of data that must contend with complications such as subjective observations, biases, and narrow observation chances while dealing with pressures such as budget, time, and the effort required [9]. In addition, the reliance on laboratory elicitable features to recognize altruism propensities hinders the progress of the field of altruism. Using such labor intensive methods for eliciting altruism essentially prevents scientists from recognizing behavioral features based on mobility or communication traces that range over time and space (e.g. day-night call ratio, diversity of locations visited) to predict one's propensities to behave altruistically.

Recently, mobile phones along with sensor-based data have been used by scientists to construct rich and individualized models of human behavior in social, spatial, and temporal settings and link them to depression, happiness evaluations, and school GPA (Grade Point Average) [10], [11], [12]. This progress motivates utilizing phone-based models for predicting altruism too. Such a phone-based method, if successful at predicting altruism propensity, may offer a low-cost, faster, scalable, and automatic process for making insights into altruism levels of billions of users in the big data era.

Hence, in this work, we systematically identify the associations between phone-based behavioral indicators and altruism and quantify the predictive power of such “personal big data” in inferring a person's propensity to behave altruistically.

- G.F. Bati is with the Computer Engineering Department, Umm Al-Qura University, P.O. Box 6151, Makkah, Saudi Arabia 21955 and the Electrical and Computer Engineering Department, Rutgers, The State University of New Jersey, Piscataway, NJ 08854. E-mail: gfbati@uqu.edu.sa.
- V.K. Singh is with the School of Communication and Information, Rutgers, The State University of New Jersey, New Brunswick, NJ 08901 and the MIT Media Lab, Cambridge, MA 02139. E-mail: v.singh@rutgers.edu.

The main contributions of this work are three-fold:

1. To motivate and ground the usage of phone-based features for inferring altruism propensities;
2. To identify the associations between long-term “in the wild” socio-mobile behavior and altruism propensities; and
3. To define a machine learning model that automatically infers a person’s propensity to be altruistic.

The rest of the paper is organized as follows. We present the related work on the issue in the following section. Then, we explain and analyze the study employed in this paper. Next, the attained results along with their implications and limitations are elaborated on. Finally, we conclude the paper and suggest some potential future work.

## 2 RELATED WORK

Mobile phones (cellphones or smartphones) are considered a crucial and prime communication device by billions of people around the globe. Due to the fact that most of current mobile phones are equipped with countless sensors, there exist in the literature many studies in various fields that use sensors on mobile phones for several purposes (mobile phone sensing) [13]. In fact, recent line of work on *phoneotypic* modeling [14] suggests that a collection of phone-based behavioral features can create a unique signature for an individual that can predict other facets of the individual’s life. Also, some researchers reflect on mobile phones to be a “vast psychological questionnaire that we are constantly filling out, both consciously and unconsciously” [15]. This was one of the motivations for us to collect the dataset of this work and study its associations with altruistic propensities.

Altruism is an important social concept and has been studied across multiple disciplines (e.g., sociology, psychology) [6], [7], [8]. Given its importance and interconnections with human behavior in technology mediated environments, it has also started receiving some attention in the computational and mobile computing literature. For example, in [16], the authors used attachment transfer theory to understand reciprocal altruism for tourism online shopping using mobile phones. The impact of altruism, topologies, and traffic patterns on mobile social networks have been studied and modeled in [17]. In [18], the authors studied altruism in a delay tolerant network (DTN) based mobile social network application. And finally, in [19] the authors have argued the case for explicitly modeling altruism levels of individuals in peer-to-peer Internet Streaming Broadcast applications.

However, there are, as yet, no efforts that utilize phone-based data to create automated machine-learning models for individual altruistic propensities. This is surprising; given multiple recent studies that have suggested that mobile phone-based features may indeed be predictive of multiple behaviors and behavioral propensities for human beings (e.g., [13], [14], [20]). For instance, the authors in [14] have reported

associations between individual propensities to cooperate and phone behavior. However, altruism and cooperation are two essentially different sociological concepts. Although altruism assumes a cost for the benefactor and advantage to the beneficiary, cooperation merely predicates benefit to the beneficiary, the benefactor might also benefit from the transaction [21]. As a consequence, a technical protocol inspired by [19] which considers cooperation would need to keep a ledger of favors given and received between agents (tit for tat behavior), while the one focusing on altruism would simply need to quantify an individual agent’s desire to help others.

A preliminary (6 page) version of this paper appeared in the conference proceedings of the 3rd IEEE International Conference on Internet of People (IoP 2017; co-located with UIC 2017) [22]. The current 11 page version expands upon it in multiple ways. We have considered and studied more *phoneotypic* features (24 instead of 14) in this version. Next, we have developed and tested multiple regression approaches for modeling numeric altruism scores, which were not part of the previous version. Also, the previous version focused on automatic prediction of the altruism class based on a split around the median. As such a median splitting method has limitations in terms of its ability to capture the underlying dynamics of the data due to the arbitrary split point selection. In other words, median split is variable-oriented, not people-oriented [23], [24], [25]. Hence, we now use “k-Means++” [26] (an unsupervised machine learning algorithm) to group the participants into naturally occurring clusters/categories based on their altruism score before evaluating phone data-based models to infer the right altruism category for individuals.

## 3 STUDY

We study the interconnections between altruism and phone-based features on the data collected as part of the Rutgers Well-being Study. This study was a ten-week field and lab study conducted in Spring 2015 including 55 participants, most of whom were undergraduate students from Rutgers, The State University of New Jersey.

Initially, all participants were invited to sign a consent agreement to participate in the study and install an Android mobile app. The mobile app could record their call, SMS, and GPS logs, not content for privacy and ethical considerations. The participants were requested to be present in-person for three sessions where they filled out a number of surveys concerning their health, well-being, altruism, and some demographics. There was a compensation of \$20, \$30, and \$50 respectively for attending the three sessions. An approach of increasing the compensation for each session was adopted in an effort to reduce the dropout rate over the ten-week period of the study. We use here the altruism and demographics surveys for their relevance to this work.

Participants’ privacy was an utmost priority; hence, anonymized IMEI numbers were used to recognize the par-

ticipants. Also, the dataset was hashed before analysis. The participation in the study was voluntary and the participants could withdraw from the study at any time. All staff who handled the data in this study were trained and certified in human subject research.

The participants' ages varied from 18 to 21 years. Of these, 35 were men and 20 were women. Most of the participants were single and the median of their families' income ranges between US \$50,000 to \$74,999. Altruism propensities (tendencies) were quantified using a survey (details follow), whereas the phone-based features have been attained from an app installed in their Android mobile phones. The app was developed using the "Funf in a box" framework [27] and was released via a URL shared with the study participants. We decided not to upload this app to Google Play Store to make sure that no one beside the participants has an access to the app or its data. Fig. 1 shows a screenshot of the app. The goal of this study is to test the feasibility of eventually replacing such in-lab surveys with automated phone-based methods. Hence, as a first step in this direction, here we test the associations between phone data and such survey scores.

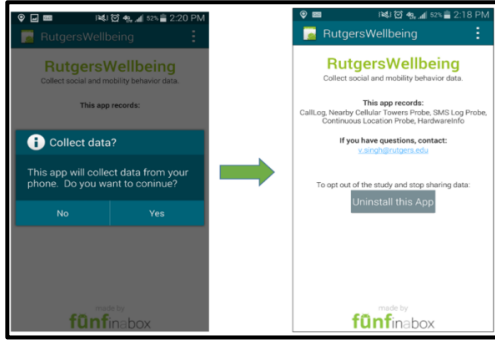


Fig. 1. Screenshot of the Android App.

The 55 participants made a total of 28,132 calls with an average of about 511 and a median of 312 calls per participant and exchanged 187,720 SMS messages with an average of 3,413 approximately and a median of 2,423 per participant, and visited 14,905 unique locations with an average of 271 and a median of 284 per participant during the period of the study (10 weeks). TABLE 1 gives a summary of the total, mean, and median for calls, SMS, and locations.

TABLE 1  
SUMMARY OF CALLS, SMS, AND LOCATIONS IN THIS STUDY

Feature	Total	Mean	Median
Calls	28,132	511	312
SMS	187,720	3,413	2,423
Unique Locations	14,905	271	284

### 3.1 Altruism Descriptor

The literature has several ways of quantifying an individual's altruism propensity. For example, games in controlled lab settings like public goods game (PGG), Trust Game, and

Dictator Game represent one way of quantifying altruism [28], [29]. Also, surveys that draw individual's behavior in prepared scenarios present another way. Additionally, another way is a combination of game experiments and lab surveys [30]. In this paper, we decided to use a well-known survey to measure altruism: "The Self-report Altruism Scale" (SRA) by Rushton et.al [31]. The survey has 20 questions whose responses scaled from (5) "very often" to (1) "never" on a five points scale. Examples of the questions are: "I have given money to a charity" and "I have pointed out a clerk's error (in a bank, at the supermarket) in undercharging me for an item" [31]. Besides the widespread adoption of the survey (over 800 citations as per Google Scholar), we chose this survey as the nature of these questions is not restricted to a specific context and the results could be interpreted in a wide variety of everyday applications. Also, SRA has an adequate validity correlations with related measures and a high reliability of  $\alpha = 0.80$  [32].

Since the survey has 20 questions worth 5 points each, the maximum theoretical altruism score is 100. In the considered sample, the maximum is found to be 95, the minimum is 31, the median is 50, and the mean is 54.07 as shown in TABLE 2.

TABLE 2  
SUMMARY OF ALTRUISM SCORES

Minimum	Maximum	Median	Mean
31	95	50	54.07

### 3.2 Demographic Descriptors

The participants were surveyed about their demography. We collected the following information: age, gender, marital status, race, level of education (school), and level of family's income.

### 3.3 Mobile Phone Data Features

To come up with a good representation of an individual's social-mobile behavior, we surveyed the related literature which focuses on connecting phone behavior with individual behaviors and social outcomes (e.g., [11], [13], [16], [19]). For example, social capital as a concept is connected with both phone use behavior [33] and altruism [34]. Social capital often comes in two variants: bridging and bonding [35]. Hence, we link the concepts of weak and strong ties to bridging and bonding social capital to predict one's propensity to altruism [14], [36], [37]. We use Call and SMS logs to represent the features that carry "social traits" concepts for mobility and altruism and their interconnections [38]. We, also, consider location logs (physical movements) as a proxy of social behavior for it has been used previously to comprehend human social behaviors [39] and human social and geo-spatial behavior are inherently connected with each other [40].

Based on the Call, SMS, Location data collected from the app, we define the following set of *phoneotypic* features (i.e. phone-based features) (N=24) as presented in TABLE 3:

TABLE 3  
SUMMARY OF PHONEOTYPIC FEATURES IN THIS STUDY

Feature	Definition
Level of Social Activity	Social Activity (Call, SMS, Location) = $\sum$ Activity Total Call Duration = $\sum$ Time Spent on I/O calls
Diversity	Diversity (Call, SMS, Location): $D_i = -\sum_j p_{ij} \log_b p_{ij}$
Novelty	Novelty (Call, SMS, Location): Percent New Contacts = $\frac{\sum \text{New Contacts}}{\sum \text{All Contacts}} \times 100$
Reciprocity	IOR = $\frac{\text{Incoming communication count}}{\text{Outgoing communication count}}$ Missed Call Percentage = $\frac{\sum \text{missed calls}}{\sum \text{calls}} \times 100$
Strong and Weak Ties Engagement Ratio	STR = $\frac{\sum \text{communication for highest } (\frac{1}{3}) \text{ contacts}}{\sum \text{communication}} \times 100$ WTR = $\frac{\sum \text{communication for lowest } (\frac{1}{3}) \text{ contacts}}{\sum \text{communication}} \times 100$
Temporal Rhythms	$\frac{\sum (\text{Call, SMS, Location}) \text{ when productive (8am to 8pm)}}{\sum (\text{Call, SMS, Location}) \text{ when relaxed (8pm to 8am)}}$ $\frac{\sum (\text{Call, SMS}) \text{ in weekdays}}{\sum (\text{Call, SMS}) \text{ in weekends}}$

### 3.3.1 Level of Social Activity

Level of Social Activity represents the activity of a user as obtained through counting exchanged phone calls, messages, and unique visited locations. A higher count of social activity level suggests an active user. The visited locations were updated hourly to balance between getting an idea about the pattern of a user's movement and their phones' battery life. To avoid getting the same amount of locations per participant (24 locations/day), we focus on unique locations. The location data were gained from a mobile phone's GPS as <latitude, longitude> tuple at fourth decimal point resolution, which roughly corresponds to 10m by 10m blocks [14], [41]. We are not only considering the total amount of calls, but also the total durations of such calls because they are related to social activity. We assume that a person who makes or receives (I/O) numerous long calls may have more social life and this may be associated with being more altruistic [6]. Thus, we consider the following features:

$$\begin{aligned} \text{Social Activity (Call, SMS, Location)} &= \sum \text{Activity} \\ \text{Total Call Duration} &= \sum \text{Time Spent on I/O calls} \end{aligned}$$

### 3.3.2 Diversity (Calls, SMS, Location)

We are not merely considering quantifying calls, SMS messages, unique locations, but also the diversity (measured as Shannon Entropy) for each one of them, as such a diversity metric has been reported to be associated with various personal well-being outcomes and personality traits [42], [43].

Diversity (Call, SMS, Location):

$$D_i = -\sum_j p_{ij} \log_b p_{ij}$$

Where  $p_{ij}$  is the percentage of social events involving individual 'i' and contact 'j', and 'b' is the total number of such contacts.

### 3.3.3 Novelty (Call, SMS, Location)

The growth of networks plays an important role in social capital [44]. Hence, we, also, consider "new contacts" that are not present in the first four weeks of the data collection period. This feature quantifies how much time users devote to their new contacts as compared to their frequent contacts.

Novelty (Call, SMS, Location):

$$\text{Percent New Contacts} = \frac{\sum \text{New Contacts}}{\sum \text{All Contacts}} \times 100$$

### 3.3.4 Reciprocity (Call, SMS)

Besides the frequency of communication, the ease with which communication is conducted is also an important property of an individual's social behavior. We anticipate approachability of individuals to be associated with their civic participation and social capital levels [33]. Such social capital levels have been associated with altruism [34]. Hence, we compute the ratio of incoming to outgoing calls and SMS text messages and also the percentage of missed calls as follows.

In Out Ratio (Call, SMS):

$$\text{IOR} = \frac{\text{Incoming communication count}}{\text{Outgoing communication count}}$$

$$\text{Missed Call Percentage} = \frac{\sum \text{missed calls}}{\sum \text{issued}} \times 100$$

### 3.3.5 Strong and Weak Ties Engagement Ratio (Call, SMS, Location)

It is anticipated that a person would devote at least 33% of their time with their top third contacts. Nevertheless, a higher score like 80% may indicate an individual's preference to pointedly engage more with strong ties rather than spreading the communication effort more equally among all ties. We were inspired by prior studies linking strength of ties and altruism [45], [46], and conjecture that the relative spread (or concentration) of communication with such strong ties may be a predictor of one's propensity to be altruistic.

$$\text{STR} = \frac{\sum \text{communication for highest } (\frac{1}{3}) \text{ contacts}}{\sum \text{communication}} \times 100$$

$$\text{WTR} = \frac{\sum \text{communication for lowest } (\frac{1}{3}) \text{ contacts}}{\sum \text{communication}} \times 100$$

### 3.3.6 Temporal Rhythms (Call, SMS, Location)

Prior literature has connected animal rhythms and circadian cycles and altruism [47]. The characterization of different individual's chronotype - the tendency for the individual to sleep at a particular time during a 24-hour period - colloquially "morningness" or "eveningness" has been connected with cheating and Machiavellianism [48]. When we asked some of the participants (mainly students) about their daily activities, times when they become productive, and times when they tend to play or sleep (relax), we found that there are two main states: "productive" state from 8 am to 8 pm; "relax" state from 8 pm to 8 am. To get more in-

sights out of these features, we added another layer of characterization for these states.

$$\frac{\sum(\text{Call, SMS, Location})_{\text{when productive(8am to 8pm)}}}{\sum(\text{Call, SMS, Location})_{\text{when relaxed(8pm to 8am)}}$$

We added another layer of characterization for the abovementioned two states of the daily activity ratio (productive and relaxed) to get more insights out of these circadian rhythms by quantifying the weekdays (Monday to Friday) to weekends (Saturday and Sunday) communication (Call, SMS) ratio.

$$\frac{\sum(\text{Call, SMS})_{\text{in weekdays}}}{\sum(\text{Call, SMS})_{\text{in weekends}}}$$

## 4 RESULTS AND DISCUSSION

Multiple applications vary in their requirements of either estimating an exact numeric for altruism score (e.g. for studying altruism levels in social science studies) or working with broader classifications of altruism score (e.g. for suggesting different default preferences for bandwidth sharing). Hence, we consider both types of applications by undertaking linear regression and classification analyses as follows.

### 4.1 Building a Regression Model for Altruism

We first consider predicting altruism level as a regression problem; that is, predicting an outcome variable (i.e., altruism level) from a set of input predictors (i.e., phone-based features). We use the LASSO (Least Absolute Shrinkage and Selection Operator) regression approach to undertake this process [49]. LASSO is a specialized form of regression suitable for scenarios where there are relatively more number of features for a given sample size. It tries to minimize overfitting by penalizing the presence of many features in the final model. It has been applied in similar contexts (in terms of sample size, number of features, and application) in recent human-centered/ubiquitous computing research [12], [33]. Similarly, following [12], [33], we assess the regression models using the metrics of correlation scores (between predicted and actual outcome variables), the Root Mean Square Error (RMSE), and the Mean Absolute Error (MAE). While a higher correlation (closer to 1) suggests a higher predictive ability of the considered models, smaller RMSE and MAE are preferred as they show that the predictions are closer to the ground truth altruism survey.

We ran and tested three different regression models: one with the demographic features only, another one with the *phoneotypic* (phone-based) features only, and a third one with a combination of both types of features. All demography features were found to be significant (Age, Gender, School, Race, and Income) in the demography only model except “Marital”. All *phoneotypic* features were found to be significant (except Weak Ties (Location)) in the *phoneotype* only model. Finally, all demography and *phoneotype* features were significant (except In Out ratio (Call), Weak Ties (Location), and Race) in the combined model. The imple-

mentation was undertaken using R 3.4.1 [50] and its Lars 1.2 package [51]. TABLE 4 presents the results for the evaluation in terms of the three metrics considered.

The demography based model obtained a correlation of 0.33 between the predicted and actual altruism values and RMSE and MAE scores of 14.69 and 11.67 respectively. The low - but significant - scores for the “demography only” model indicate that the demographic features can explain some (but not a lot) of variance in the altruism levels. Phone-based model performed much better with a correlation score of 0.75 (RMSE=10.26; MAE=7.99).

The combined model using *phoneotype* and demography features performed the best in terms of all three metrics and the predicted altruism was found to have 0.81 correlation with the actual altruism scores (RMSE=9.18; MAE=7.24). An MAE of 7.24 implies that the predictions are within  $\pm 7.24$  of the absolute value of the altruism scores obtained by the survey (ground truth). Since the altruism scores obtained by the survey vary between 31 and 95 as shown in TABLE 2, ranges of  $\pm 7.24$  can reflect a quite reasonable approximation.

Also, we see that the *phoneotype* model and “*phoneotype* + demography” (both) models yield considerably better models than the demography-based model. However, the demographic features were useful in increasing the correlation score for the *phoneotypic* model, thus suggesting that *phoneotypic* features and demography features are not merely proxies for each other, but rather add newer information when combined.

TABLE 4  
MODELING ALTRUISM USING DIFFERENT REGRESSION MODELS

Model Type	Correlation	RMSE	MAE
Demography Only	0.33	14.69	11.67
<i>Phoneotype</i> Only	0.75	10.26	7.99
Both	0.81	9.18	7.24

### 4.2 Building a Classification Model for Altruism

We aim to build and test a classification model capable of predicting altruistic propensities. The literature suggests various methods to cluster (group or categorize) such data, including standard median splits and extreme group analysis [24]. Standard median splits dichotomize continuous variables into two groups: “low” which is lower than the median value (50) of the data and “high” which is greater than the median as we have done in our previous work [22]. However, such median splits have limitations in the sense they are often unable to capture the underlying dynamics of the observed phenomena because they are variable-oriented and not people-oriented [23], [24], [25].

Hence, in this version, we are going to use “unsupervised machine learning” to cluster the participants into naturally occurring groups based on their altruism scores. Specifically, we use k-Means++ [26] clustering algorithm to find the optimal clusters. The algorithm is initialized by

choosing the first center randomly. Then, the succeeding centers are selected from the remaining points based on the squared distance from the closest center. We ran the algorithm ten times with 300 maximum iterations per each algorithm run. An important consideration for k-Means++ algorithm is the choice of the number of clusters (k) to be used by the algorithm. Literature suggests multiple methods including: Silhouette scores, Bayesian Information Criteria (BIC), and “elbow method” for identifying the right number of clusters [52].

Here, we considered two different methods (Silhouette scores and BIC) for this process. Silhouette scores (higher score is better) compare the average distance to elements in the same cluster with the average distance to elements in other clusters [53]. We implemented this procedure in R 3.4.1 [50] and its package ClusterR 1.0.6 [54]. We found that the best k equals two as shown in Fig. 2. This method generated two clusters, one of which contains the altruism scores 31 to 58 (N=36 participants), second of which contains the rest of the scores (60 to 95) (N=19 participants). For the ease of interpretation, we refer to these groups as “altruism group A” and “altruism group B” respectively.

Using the BIC criteria (lower score is better) with the same ClusterR package to identify the optimal number of clusters for k-Means++, however, suggested the optimal number of cluster to be three. (Please refer to Fig. 3.) The first identified cluster contains the altruism scores from 31 to 45 (N=17 participants), second cluster contains the altruism scores 46 to 62 (N=24 participants), and the third cluster contains the rest of the scores (65 to 95) (N=14 participants).

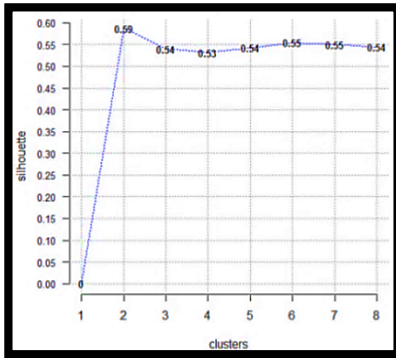


Fig. 2. Optimal Number of Clusters for k-Means++ (Silhouette Score).

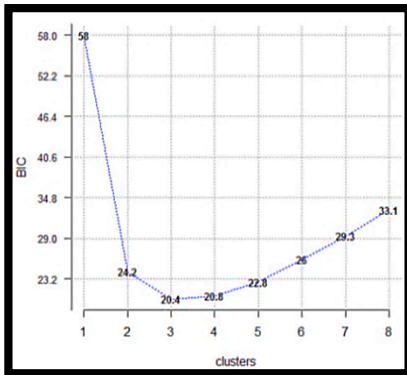


Fig. 3. Optimal Number of Clusters for k-Means++ (BIC).

We used Orange3-3.6.0 [53], [55] to build the models which could automatically identify the altruism group (e.g. “altruism group A” or “altruism group B”), which an individual belongs to. We built 3 types of models based on k=2: one with the demographic features only, another one with the *phoneotypic* features only, and a third one with a combination of both types of features. We used information gain (reduction of entropy) [53] to rank the best subset of the (24 *phoneotype* + 6 demography) features described in the preceding section. For optimal feature subset selection, we first ranked all the features based on information gain [53]. Then, we considered models of up to ten features (a third of the available pool) wherein each model was the collection of top “j” features and j is in the range (1, 10). The optimal subset was the one with highest performance amongst the considered models. The resulting feature sets in each of the cases is shown in TABLE 5.

The abovementioned features were used to test out several well-known machine learning algorithms for classification. Specifically, we used Naïve Bayes, Random Forest, CN2 Rule Induction, Logistic Regression, and kNN (k-Nearest Neighbors) with a leave-one-out cross validation method to balance between the learning opportunities and the generalizability of results from the data. Also, we used Zero-R (Constant or Majority) without any cross validation which simply classifies all the instances into the majority class as a baseline to facilitate interpreting the results. TABLE 6 offers a comparison of the results. It is worth noting that AUC stands for (Area Under the Receiver Operating Characteristic Curve) and CA means (Curve Accuracy) [53]. Moreover, F1 represents the harmonic mean between precision and recall [56]. A higher score (closer to 1) is better in each case and the Zero-R scores give a sense of the baseline expected performance.

TABLE 6 shows that the demography-based model returned the best curve accuracy (CA) of 70.9%, top AUC of 0.656, and F1 score of 0.486. The *phoneotype*-based model generated a better model than the demography model whose best accuracy is 78.2%, AUC is 0.756, and F1 is 0.667. The best results however were obtained by a combined model (demography + *phoneotype* data) that yielded an accuracy of 78.2%, AUC of 0.798, and F1 of 0.667.

TABLE 5

FEATURES SELECTED FOR VARIOUS PREDICTION MODELS (K=2)

<b>Demography Only</b>	Income, School (Level of Education)
<b>Phoneotype Only</b>	Total Call Duration, Missed Call Percentage, Number of New Contacts (SMS)
<b>Both</b>	Total Call Duration, Income, Missed Call Percentage, Number of New Contacts (SMS), School

We repeated a similar process of building three predictive models based on k=3. We found that the best results were obtained when we selected the following features in each of the models as shown in TABLE 7. TABLE 8 shows that the demography-based model returned the greatest accuracy (CA) of 45.5%, top AUC of 0.577, and F1 score of

0.517. The *phoneotype*-based model generated a better model than the demography model whose best accuracy is 56.4%, AUC is 0.777, and F1 is 0.717. The best results were obtained by a combined model (demography + *phoneotype* data) that yielded an accuracy of 65.5%, AUC of 0.816, and F1 of 0.760.

Note that these results in terms of accuracy are lower for the three-way classification problem compared to the two-way classification problem. However, a three-way classification problem is in general a harder problem than two way classification, and the much lower baseline (Zero-R) scores may help interpret the performance gain obtained by the phone-based models.

TABLE 7  
FEATURES SELECTED FOR VARIOUS PREDICTION MODELS (K=3)

<b>Demography Only</b>	Income, School (Level of Education), Age
<b>Phoneotype Only</b>	Missed Call Percentage, Weak Ties (SMS), Weekday Weekend SMS Ratio, Number of New Contacts (SMS), Weak Ties (Call)
<b>Both</b>	Income, Missed Call Percentage, Weak Ties (SMS), School, Age, Weekday Weekend SMS Ratio, Number of New Contacts (SMS), Weak Ties (Call)

From the aforementioned results, we can clearly observe that *phoneotypic* features considerably outperform demography-based ones in prediction's accuracy corroborating the findings from the regression analysis. Further, similar to the regression analysis, the *phoneotypic* (i.e., phone-based) behavioral features were not merely a replacement for demographic features, as the combined models yielded higher performance as compared to the individual models.

It is also clear that the *phoneotypic* model outperformed the baseline Zero-R model. In the case of  $k=2$ , the *phoneotypic* model performed 59.6% better than the baseline model in terms of AUC, 19.4% better in terms of accuracy, and 28.8% better in terms of F1 score. In the case of  $k=3$ , the *phoneotypic* model performed 63.2% better than the baseline model in terms of AUC, 50.2% better in terms of accuracy, and 186.8% better in terms of F1 score.

Hence, we note that a phone-features based model beats baseline majority classification and also goes beyond static demographic descriptors (e.g. age, gender, education, income) for predicting altruism propensities. This underscores the potential for using phone-based (or *phoneotypic*) features to build automatic classifiers for individual altruism propensities. These findings provide clear evidence of interconnections between the mobile features and altruistic propensities and also motivating further work in this direction. In effect, these results pave way for "personal big data" to expand the understanding of the associations between socio-mobile behavioral data and altruism propensities.

### 4.3 Features Associated with Altruism Propensity

Besides identifying automated methods for identifying an individual's altruism propensity levels, one of the goals of this work is to understand the socio-mobile behavior of individuals with different propensities to be altruistic. Thus, we undertook a *post-hoc* Pearson's correlation analysis using IBM SPSS 24 between the altruism scores obtained from the survey and the *phoneotypic* features. Note that the correlation analysis undertaken here is post-hoc and intended to help interpret the observed predictions, as opposed to being prescriptive in its own right.

In the interest of space, we only report the correlations that were found to be (at least marginally i.e.  $p < 0.10$ ) significant in TABLE 9. We note that people who have high altruism propensity tend to be more socially active, yet have different usage patterns for different communication modalities.

First, individuals with higher altruism propensity tend to call more often ( $r = +0.356$ ). This can be understood as altruism propensity being associated with healthy social relationships and higher call activity captures such behavior [6]. We see a similar trend in terms of new call contacts ( $r = +0.305$ ). This underscores the importance of constantly renewing and broadening one's social contacts and its associations with altruism. From a methodological perspective the dynamics of social contacts captured by this feature underscore the value of temporal features, which cannot be captured in one-time lab studies studying the phenomena of altruism.

Next, we notice that individuals with higher altruism propensity show a marked preference for engaging in phone calls with their "strong ties" as opposed to spending it equitably with all contacts ( $r = +0.252$ ). Conversely, they spend less time on calls with their "weak ties" ( $r = -0.232$ ). However, the patterns of SMS communication seem to be quite different from phone call-based communication. SMS interactions with "weak ties" were found to be positively associated ( $r = +0.260$ ) with altruism. From a methodological perspective, these results suggest the value of different modalities of data to triangulate and predict human traits. From a conceptual perspective, these observations corroborate previous studies which suggest that altruism is a pro-social trait associated with higher social capital including both its bridging and bonding variants [45], [46].

We note that the correlation scores found are relatively modest and unlikely to be prescriptive in their own right. At the same time, the correlation scores for each of the significant features in TABLE 9 were also found significant in LASSO regression as explained earlier. Hence, this allows for triangulation across methods to interpret the results observed.

TABLE 6  
PREDICTION RESULTS FOR ALTRUISM LEVELS USING DIFFERENT ALGORITHMS (K=2)

Method	Demography Only			Phoneotype Only			Both		
	AUC	CA	F1	AUC	CA	F1	AUC	CA	F1
NaïveBayes	0.594	0.655	<b>0.486</b>	<b>0.756</b>	0.727	0.595	<b>0.798</b>	<b>0.782</b>	<b>0.667</b>
RandomForest	0.594	0.673	0.471	0.664	0.636	0.375	0.679	0.691	0.452
CN2 Rule	<b>0.656</b>	<b>0.709</b>	0.429	0.508	0.618	0.276	0.587	0.582	0.343
LogRegression	0.558	0.582	0.303	0.586	0.618	0.160	0.598	0.600	0.154
kNN	0.459	0.655	0.000	0.719	<b>0.782</b>	<b>0.667</b>	0.719	<b>0.782</b>	<b>0.667</b>
Zero-R	0.500	0.655	0.518	0.500	0.655	0.518	0.500	0.655	0.518

TABLE 8  
PREDICTION RESULTS FOR ALTRUISM LEVELS USING DIFFERENT ALGORITHMS (K=3)

Method	Demography Only			Phoneotype Only			Both		
	AUC	CA	F1	AUC	CA	F1	AUC	CA	F1
NaïveBayes	0.535	0.273	0.279	<b>0.777</b>	<b>0.564</b>	<b>0.717</b>	<b>0.816</b>	<b>0.655</b>	<b>0.760</b>
RandomForest	0.519	0.382	0.431	0.569	0.455	0.538	0.558	0.364	0.464
CN2 Rule	<b>0.577</b>	<b>0.455</b>	0.435	0.624	0.436	0.576	0.639	0.473	0.571
LogRegression	0.526	0.436	<b>0.517</b>	0.414	0.255	0.393	0.528	0.491	0.566
kNN	0.507	0.418	0.412	0.593	0.418	0.500	0.606	0.436	0.528
Zero-R	0.500	0.436	0.265	0.500	0.436	0.265	0.500	0.436	0.265

TABLE 9  
PEARSON'S CORRELATION FOR ALTRUISM (\*\* 0.01, \* 0.05, ° 0.10)

Phoneotypic Feature	Pearson's Correlation	Significance (P-value)
Social Activity Level (Call)	0.356	0.008**
Number of New Contacts (Call)	0.305	0.024*
Strong Ties (Call)	0.252	0.064°
Weak Ties (Call)	-0.232	0.089°
Weak Ties (SMS)	0.260	0.055°

#### 4.4 Discussion

The three forms of analysis implemented in this work (regression analysis, classification, and Pearson's correlation analysis) suggest that machine learning, and data analytics approaches generally, can be utilized to infer individual altruism propensity based on phone metadata to a substantial extent. The regression analysis can estimate the individual altruism propensities with high correlation (0.75), and within a margin of  $\pm 7.99$  over a range of 31 to 95. Accompanying phone features with demographic data, where available, could yield even better performance. For instance, the classification analysis yielded up to about 80% accuracy (AUC=0.782; F1=0.667) based on such combination of *phoneotype* and demography features.

Given the small sample size, we focus here on exploring general patterns and trends over the three analysis techniques (regression, classification, and correlation). We can observe a consistency in the results across the three analysis types as well as the two variants (k=2 and k=3) for classification, suggesting that socio-mobile signals as observed via a phone (*phoneotype*)

could indeed be used to infer altruism propensity of an individual. The results contribute to the growing literature on using "personal big data" to characterize multiple traits of human beings [12], [14], [27], [33], [42]. At the same time, they motivate further work to study altruism propensity using socio-mobile behavioral data.

##### 4.4.1 Privacy of User Data and Ethical Considerations

To insure and maintain the privacy of the participants, we followed the best practices in the human subjects' research that require hashing and anonymizing all data before analysis. Also, no one from the research team under any circumstance had an access to private data like the exact phone number of a participant or the content of the calls or SMS messages. The Android app collecting the data requires lesser permissions than many of the popular apps available at Google Play Store (e.g. WhatsApp, Instagram).

We also note the moral and ethical considerations in giving a person a score based on their propensity to be altruistic. History repeats itself; similar reservations have been raised up about the conventional paper survey approaches with a similar objective, and likewise automatic systems which use social media and phone data to give health, well-being, or similar scores to people [57]. Rather than waiting for perfect privacy and ethics guidelines to emerge around these topics, we posit that studies like ours can help broaden the understanding around the prospects of using "personal big data" to create personalized sociological profiles of individuals and inform the discussion in the research community around them [58].

#### 4.4.2 Limitations

This study has three limitations: 1. homogeneity of the sample (participants were mostly undergraduate students from the same institution), 2. small sample size (55) while having large number (24) of potentially collinear features in regression analysis, and 3. Inability to establish causality. Bearing in mind these limitations, we will be cautious in generalizing the findings obtained until they are verified at scale over representative sample populations. To overcome these limitations, we used LASSO regression which deals with such situations of having relatively more number of features for a sample trying to minimize overfitting by penalizing the use of many features [49]. Furthermore, we plan to repeat the study in the future considering a larger and more diverse sample.

Despite these limitations, to the best of our knowledge, this is the first line of work to analyze the links between altruism levels and phone-based socio-mobile behavior (*phoneotype*). The obtained results in this first of its kind effort are thus encouraging, and have demonstrated the potential of “personal big data” for predicting altruism levels of individuals.

#### 4.4.3 Implications

The results open the doors to a methodology that, with refinements and validation, could be used at scale. Smartphones are now actively used by more than 1.4 billion users, and hence the proposed method could potentially be applied to estimate the altruistic levels for billions of individuals. In this sense, this work comes to an agreement with the idea highlighted in the smartphone psychology manifesto, which states that “... smartphones could transform psychology even more profoundly than PCs and brain imaging did” ([59], p.1).

In future, this work could also have multiple implications for social scientists, economists, mobile phone service providers, and policy designers. For example, the suggested methodology could help social scientists study altruism at scale in the society. Besides identifying connections between spatial and temporal behavior and altruism, a scalable methodology to study altruism could allow for asking questions regarding the spread of altruism in networks of billions of individuals, which are simply not possible with current survey or lab-based methods.

Similarly, as mobile phones are increasingly used, both, as user end-points and as mediators of technology, modeling a person’s altruistic propensities automatically could be helpful in supporting various socio-technical applications under the Internet of People vision [60]. Such an Internet-of-People vision explicitly requires the creation of a “sociological profile” [60] for the participants and the proposed method for inferring altruistic propensities could be used, for example, to identify a person’s default preferences in peer-to-peer networking, file sharing, or human-computation based tasks. Further, with multiple bots negotiating services and conditions for users in the emerging social Internet of Things scenarios, having such a sociological profile could be useful to suggest default settings in multiple scenarios, from something as simple as setting the right room temperature in shared workspaces to the default “tipping” amount in dinner payments.

## 5 CONCLUSIONS AND FUTURE WORK

We have proposed a new “personal big data” based methodology to predict individual’s propensities for altruism using phone features as an alternative to classical methods like surveys and lab experiments. Using these features allowed us to build prediction models by means of machine learning classification algorithms whose AUC, accuracy, and F1 are promising and encouraging. Specifically, the best performing model for the considered two-class classification problem yields AUC = 0.798, accuracy = 78.2%, and F1 = 0.667 and the corresponding scores for the three-class problem are AUC = 0.816, accuracy = 65.5%, and F1 = 0.760. To the best of our knowledge, there has been no earlier line of work that analyzes the associations between individual altruism propensity levels and phone-based behavioral features. Thus, these results pave way to further research on leveraging ubiquitous sensing data for automatic altruism inference with applications in security, networking, business, and well-being.

While we consider the results to be early and exploratory, the proposed approach can be enhanced in future work by including a larger number of participants, more detailed phone-based features, and considering larger time durations. Furthermore, the suggested phone-based method could be expanded to study and predict other personal behaviors and traits such as trust, gratitude, compassion, and happiness. Taken together such methods open ways to better model human beings based on ubiquitous sensing and act as a building block towards a healthier and happier society.

## Acknowledgment

G.F.B. would like to thank Umm Al-Qura University, Makkah, Saudi Arabia, for partially sponsoring this work and providing the fellowship to pursue his graduate studies at Rutgers, The State University of New Jersey. The Authors would like to thank Cecilia Gal, Padmapriya Subramanian, Ariana Blake, Suril Dalal, Sneha Dasari, Isha Ghosh, and Christin Jose for assisting in conducting the study and processing the data.

## References

- [1] C. Gurrin, A. F. Smeaton and A. R. Doherty, “Lifelogging: Personal big data,” *Foundations and Trends in Information Retrieval*, vol. 8, no. 1, pp. 1-125, 2014.
- [2] A. Ali, J. Qadir, R. u. Rasool, A. Sathiaselan, A. Zwitter and J. Crowcroft, “Big data for development: applications and techniques,” *Big Data Analytics*, vol. 1, no. 2, 2016.
- [3] V. D. Blondel, A. Decuyper and G. Krings, “A survey of results on mobile phone datasets analysis,” *EPJ Data Science*, vol. 4, no. 10, 2015.
- [4] K. G. Tarakji, O. M. Wazni, T. Callahan, M. Kanj, A. H. Hakim, K. Wolski and B. D. Lindsay, “Using a novel wireless system for monitoring patients after the atrial fibrillation ablation procedure: the iTransmit study,” *Heart Rhythm*, vol. 12, no. 3, pp. 554-559, 2015.
- [5] A. Decuyper, “On the research for big data uses for public good purposes Opportunities and challenges,” *Netcom*, vol. 30, no. 3-4, pp. 305-314, 2016.
- [6] The Greater Good Science Center at UC Berkeley, “What Is Altruism?,” The Greater Good Science Center at UC Berkeley, [Online]. Available: <http://greatergood.berkeley.edu/topic/altruism/definition>. [Accessed

- 04 April 2017].
- [7] S. Okasha, "Biological Altruism (Stanford Encyclopedia of Philosophy)," 21 July 2013. [Online]. Available: <https://plato.stanford.edu/archives/fall2013/entries/altruism-biological/>. [Accessed 4 April 2017].
  - [8] C. Batson, *The altruism question: Toward a social-psychological answer*, New York: Psychology Press, 2014.
  - [9] J. Giles, "Computational social science: Making the links," *Nature*, vol. 488, no. 7412, pp. 448-450, 2012.
  - [10] S. Consolvo, D. McDonald, T. Toscos, M. Chen, J. Froehlich, B. Harrison, P. Klasnja, A. LaMarca, L. LeGrand, R. Libby and I. Smith, "Activity sensing in the wild: a field trial of ubifit garden," in the SIGCHI Conference on Human Factors in Computing Systems, Florence, Italy, 2008.
  - [11] E. Shmueli, V. Singh, B. Lepri and A. Pentland, "Sensing, understanding, and shaping social behavior," *IEEE Transactions on Computational Social Systems*, vol. 1, no. 1, pp. 22-34, 2014.
  - [12] R. Wang, F. Chen, Z. Chen, T. Li, G. Harari, S. Tignor, X. Zhou, D. Ben-Zeev and A. Campbell, "StudentLife: assessing mental health, academic performance and behavioral trends of college students using smartphones," in the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing, Seattle, Washington, US, 2014.
  - [13] W. Z. Khan, Y. Xiang, M. Y. Aalsalem and Q. Arshad, "Mobile Phone Sensing Systems: A Survey," *IEEE COMMUNICATIONS SURVEYS & TUTORIALS*, vol. 15, no. 1, pp. 402 - 427, 2013.
  - [14] V. K. Singh and R. R. Agarwal, "Cooperative Phoneotypes: Exploring Phone-based Behavioral Markers of Cooperation," in The 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp '16), Heidelberg, Germany, 2016.
  - [15] H. Grassegger and M. Krogerus, "The Data That Turned the World Upside Down," *Vice Media*, 28 January 2017. [Online]. Available: [https://motherboard.vice.com/en\\_us/article/mg9vvn/how-our-likes-helped-trump-win](https://motherboard.vice.com/en_us/article/mg9vvn/how-our-likes-helped-trump-win). [Accessed 18 September 2017].
  - [16] M. J. Kim, N. Chung, C.-K. Lee and M. W. Preis, "Why do smartphone shoppers help others on websites? The effects of attachments on reciprocal altruism," *Information Development*, vol. 32, no. 4, pp. 920-936, 2016.
  - [17] P. Hui, K. Xu, V. Li, J. Crowcroft, V. Latora and P. Lio, "Selfishness, altruism and message spreading in mobile social network," in INFOCOM Workshops 2009, Rio de Janeiro, Brazil, 2009.
  - [18] S. Hameed, A. Wolf, K. Zhu and X. Fu, "Evaluation of Human Altruism Using a DTN-based Mobile Social Network Application," in 5. Workshop Digitale Soziale Netze, Braunschweig, Germany, 2012.
  - [19] Y. Chu and H. Zhang, "Considering altruism in peer-to-peer internet streaming broadcast," in 14th international workshop on Network and operating systems support for digital audio and video NOSSDAV '04, Cork, Ireland, 2004.
  - [20] N. Lane, E. Miluzzo, H. Lu, D. Peebles, T. Choudhury and A. Campbel, "A survey of mobile phone sensing," *IEEE Communications magazine*, vol. 48, no. 9, pp. 140-150, September 2010.
  - [21] S. West, A. Griffin and A. Gardner, "Social semantics: altruism, cooperation, mutualism, strong reciprocity and group selection," *Journal of Evolutionary Biology*, vol. 20, no. 2, p. 415-432, 2007.
  - [22] G. F. Bati and V. K. Singh, "Are You Altruistic? Your Mobile Phone Could Tell," in 2017 IEEE SmartWorld Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI), San Francisco, CA, 2017.
  - [23] G. Danilo, S. MacDonald and T. Archer, "Two different approaches to the affective profiles model: median splits (variable-oriented) and cluster analysis (person-oriented)," *PeerJ* 3:e1380, 2015.
  - [24] J. DeCoster, M. Gallucci and A.-M. R. Iselin, "Best Practices for Using Median Splits, Artificial Categorization, and their Continuous Alternatives," *Journal of Experimental Psychopathology*, vol. 2, no. 2, pp. 197-209, 2011.
  - [25] D. H. Kim, N. K. Seely and J.-H. Jung, "Do you prefer, Pinterest or Instagram? The role of image-sharing SNSs and self-monitoring in enhancing ad effectiveness," *Computers in Human Behavior*, vol. 70, no. C, pp. 535-543, 2017.
  - [26] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," in The eighteenth annual ACM-SIAM symposium on Discrete algorithms, New Orleans, Louisiana, 2007.
  - [27] N. Aharony, W. Pan, C. Ip, I. Khayal and A. Pentland, "Social fMRI: Investigating and shaping social mechanisms in the real world," *Pervasive and Mobile Computing*, vol. 7, no. 6, pp. 643-659, 2011.
  - [28] F. Exadaktylos, A. M. Espín and P. Brañas-Garza, "Experimental subjects are not different," *Scientific Reports*, vol. 3, 2013.
  - [29] P. Barclay, "Trustworthiness and competitive altruism can also solve the "tragedy of the commons"," *Evolution and Human Behavior*, vol. 25, no. 4, pp. 209-220, 2004.
  - [30] J. Ermisch, D. Gambetta, H. Laurie, T. Siedler and S. C. Noah Uhrig, "Measuring people's trust," *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, pp. 749-769, 2009.
  - [31] J. Rushton, R. Chrisjohn and G. Fekken, "The altruistic personality and the self-report altruism scale," *Personality and individual differences*, vol. 2, no. 4, pp. 293-302, 1981.
  - [32] S. King Jr and M. Holosko, "The development and initial validation of the empathy scale for social workers," *Research on Social Work Practice*, vol. 22, no. 2, pp. 174-185, 2012.
  - [33] V. Singh and I. Ghosh, "Inferring Individual Social Capital Automatically via Phone Logs," *Proceedings of the ACM Human Computer Interaction*, vol. 1, no. 2, 2017.
  - [34] R. Putnam, "Social capital: Measurement and consequences," *Canadian Journal of Policy Research*, vol. 2, no. 1, pp. 41-51, 2001.
  - [35] R. D. Putnam, "Bowling alone: America's declining social capital," *Journal of democracy*, vol. 6, no. 1, pp. 65-78, 1995.
  - [36] M. S. Granovetter, "The Strength of Weak Ties," *American Journal of Sociology*, vol. 78, no. 6, pp. 1360-1380, 1973.
  - [37] D. Williams, "On and Off the 'Net: Scales for Social Capital in an Online Era," *Journal of Computer-mediated Communication*, vol. 11, no. 2, pp. 593-628, 2006.
  - [38] J. Le Galliard, R. Ferrière and U. Dieckmann, "Adaptive Evolution of Social Traits: Origin, Trajectories, and Correlations of Altruism and Mobility," *The American Naturalist*, vol. 165, no. 2, pp. 206-224, 2005.
  - [39] G. M. Vazquez-Prokopec, D. Bisanzio, S. Stoddard, V. Paz-Soldan, A. Morrison, J. Elder, J. Ramirez-Paredes, E. Halsey, T. Kochel, T. Scott and U. Kitron, "Using GPS technology to quantify human mobility, dynamic contacts and infectious disease dynamics in a resource-poor urban environment," *PloS one*, vol. 8, no. 4, 2013.

- [40] S. Scellato, A. Noulas, R. Lambiotte and C. Mascolo, "Socio-spatial properties of online location-based social networks," in Fifth International AAAI Conference on Weblogs and Social Media (ICWSM 2011), Barcelona, Spain, 2011.
- [41] D. Wang, D. Pedreschi, C. Song, F. Giannotti and A.-L. Barabasi, "Human Mobility, Social Ties, and Link Prediction," in The 17th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '11), San Diego, California, 2011.
- [42] Y. A. de Montjoye, S. S. Wang, A. Pentland, D. T. T. Anh and A. Datta, "On the Trusted Use of Large-Scale Personal Data," *IEEE Data Eng. Bull.*, pp. 5-8, 2012.
- [43] V. K. Singh, L. Freeman, B. Lepri and A. P. Pentland, "Classifying Spending Behavior using Socio-Mobile Data," *Human Journal*, pp. 99-111, 2013.
- [44] N. Lin, *Social capital: A theory of social structure and action*, New York, NY: Cambridge university press, 2002.
- [45] O. Curry, S. Roberts and R. Dunbar, "Altruism in social networks: Evidence for a 'kinship premium'," *British Journal of Psychology*, vol. 104, no. 2, pp. 283-295, 2013.
- [46] R. Bapna, A. Gupta, S. Rice and A. Sundararajan, "Trust, reciprocity and the strength of social ties: An online social network based field experiment," in Workshop on Information Systems and Economics, Shanghai, China, 2011.
- [47] P. K. Jonason, A. Jonesb and M. Lyonsb, "Creatures of the night: Chronotypes and the Dark Triad traits," *Personality and Individual Differences*, vol. 55, no. 5, p. 538-541, 2013.
- [48] A. Adan and H. Almirall, "Horne & Östberg morningness-eveningness questionnaire: A reduced scale," *Personality and Individual differences*, vol. 12, no. 3, pp. 241-253, 1991.
- [49] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, no. 1, pp. 267-288, 1996.
- [50] R Core Team, *R: A Language and Environment for Statistical Computing*, Vienna, Austria: R Foundation for Statistical Computing, 2017.
- [51] T. Hastie and B. Efron, *lars: Least Angle Regression, Lasso and Forward Stagewise*, 2013.
- [52] T. M. Kodinariya and P. R. Makwana, "Review on determining number of Cluster in K-Means Clustering," *International Journal of Advance Research in Computer Science and Management Studies*, vol. 1, no. 6, pp. 90-95, 2013.
- [53] Orange Data Mining, "Orange Visual Programming Documentation Reales 3," 09 October 2017. [Online]. Available: <https://media.readthedocs.org/pdf/orange-visual-programming/latest/orange-visual-programming.pdf>. [Accessed 14 October 2017].
- [54] L. Mouselimis, "ClusterR: Gaussian Mixture Models, K-Means, Mini-Batch-Kmeans and K-Medoids Clustering," 03 August 2017. [Online]. Available: <https://cran.r-project.org/web/packages/ClusterR/>. [Accessed 11 October 2017].
- [55] J. Demšar, T. Curk, A. Erjavec, Č. Gorup, T. Hočevár, M. Milutinović, M. Možina, M. Polajnar, M. Toplak, A. Starič, M. Štajdohar, L. Umek, L. Žagar, J. Žbontar, M. Žitnik and B. Zupan, "Orange: Data Mining Toolbox in Python," *Journal of Machine Learning Research*, vol. 14, pp. 2349-2353, 2013.
- [56] A. Zheng, *Evaluating Machine Learning Models A Beginner's Guide to Key Concepts and Pitfalls*, Sebastopol, CA: O'Reilly Media, Inc., 2015.
- [57] J. Darley and B. Latane, "The unresponsive bystander: Why doesn't he help," *New York: Appleton-Century-Crofts*, pp. 276-290, 1970.
- [58] A. Shifali, J. Yttri and W. Nilsen, "Privacy and security in mobile health (mHealth) research," *Alcohol research: current reviews*, vol. 36, no. 1, p. 143-151, 2014.
- [59] G. Miller, "The smartphone psychology manifesto," *Perspectives on psychological science*, vol. 7, no. 3, pp. 221-237, 2012.
- [60] J. Miranda, N. Mäkitalo, J. Garcia-Alonso, J. Berrocal, T. Mikkonen, C. Canal and J. M. Murillo, "From the Internet of Things to the Internet of People," *IEEE Internet Computing*, vol. 19, no. 2, pp. 40-47, 2015.

**Ghassan F. Bati** holds a BSc degree in computer engineering from Umm Al-Qura University (UQU), Holy Makkah, Saudi Arabia, 2007, an MSc degree in computer engineering from George Washington University (GWU), Washington, DC, USA, 2012, a graduate certificate in computer architecture and high-performance computing from GWU, Washington, DC, USA, 2014. He is currently pursuing his PhD in computer engineering from Rutgers University, New Jersey, USA. He works as a teaching assistant for the computer engineering department at UQU. His research interests are Reality Mining, Machine Learning, Computer Engineering Education, and High-performance Computing. He is a student member of the IEEE, the IEEE Computer Society, and the ACM.

**Vivek K. Singh** is an Assistant Professor in the School of Communication and Information at Rutgers University and a Visiting Professor at MIT Media Lab. He holds a Ph.D. in Computer Science from the University of California, Irvine (2012). He was a post-doctoral researcher at the MIT Media Lab from 2012 to 2014. His work on reality mining, cyberbullying detection, and privacy management has been published in multiple leading scientific venues (*Science*, *Proceedings of the IEEE*) and has received significant media coverage (*BBC*, *New York Times*, *Wall Street Journal*). He was selected as one of the 'Emerging Leaders in Multimedia Research' by IBM Research Labs in 2009 and he won the 2013 'Big Data for Social Good' datathon organized by Telefónica, the Open Data Institute and the MIT. He is a member of the IEEE.