

A Study of Age and Gender seen through Mobile Phone Usage Patterns in Mexico

Carlos Sarraute
Grandata Labs, Argentina
charles@grandata.com

Pablo Blanc
Mathematics Dept., FCEN, UBA
pblanc@dm.uba.ar

Javier Burroni
Grandata Labs, Argentina
javier.burroni@grandata.com

Abstract—Mobile phone usage provides a wealth of information, which can be used to better understand the demographic structure of a population. In this paper we focus on the population of Mexican mobile phone users. Our first contribution is an observational study of mobile phone usage according to gender and age groups. We were able to detect significant differences in phone usage among different subgroups of the population. Our second contribution is to provide a novel methodology to predict demographic features (namely age and gender) of unlabeled users by leveraging individual calling patterns, as well as the structure of the communication graph. We provide details of the methodology and show experimental results on a real world dataset that involves millions of users.

I. INTRODUCTION

Mobile phones have become prevalent in all parts of the world, in developed as well as developing countries, and provide an unprecedented source of information on the dynamics of the population on a national scale. In particular, mobile phone usage is starting to be used to perform quantitative analysis of the demographics of the population, respect to key variables such as gender, age, level of education and socioeconomic status (for example see [1], [2]).

In this work we combine two sources of information: transaction logs from a major mobile operator in Mexico, and information on the age and gender of a subset of the population. This allows us to perform an observational study of mobile phone usage, differentiated by gender and age groups. This study is interesting in its own right, since it provides knowledge on the structure and demographics of the mobile phone market in Mexico. We can start to fill gaps in our understanding of basic demographic questions: Are inequalities between men and women, as reported by [3], reflected in mobile phone usage (in calling and texting patterns)? What are the differences in mobile phone usage between different age groups?

The second contribution of this work is to apply the knowledge on calling patterns to predict demographic features, namely to predict the age and gender of unlabeled users. We present methods that rely on individual calling patterns, and introduce a novel algorithm that exploits the structure of the social graph (induced by communications), in order to improve the accuracy of our predictions.

Being able to understand and predict demographic features such as age and gender has numerous applications, from

market research and segmentation to the possibility of targeted campaigns (such as health campaigns for women [4]).

The remainder of the paper is organized as follows: section II provides an overview of the datasets that we used in this study. Section III describes the observations that we gathered, the insights gained from data analysis, and the differences that could be seen in CDR features between genders and age groups. In particular, very clear correlations have been observed in the links between users according to their age. In section IV we present the models that we used to identify the age and gender of unlabeled users. We show the experimental results obtained using classical Machine Learning techniques based on individual attributes, both for gender and age. We introduce a novel algorithm that leverages the links between users both in its pure graph based form (section IV-D), and combined form (section IV-E). The results of our experiments show that the pure graph based algorithm has the best predictive power. Section V concludes the paper with ideas for future work.

II. DATASET DESCRIPTION

The dataset used for this study consists of cell phone call and SMS (Short Message Service) records collected in Mexico for a period of M months ($M = 3$) by a large mobile phone operator. The dataset is anonymized. For our purposes, each CDR (Call Detail Record) is represented as a tuple $\langle x, y, t, dur, d, l \rangle$, where x and y are the encrypted phone numbers of the caller and the callee, t is the date and time of the call, dur is the duration of the call, d is the direction of the call (incoming or outgoing, with respect to the mobile operator client), and l is the location of the tower that routed the communication. Similarly, each SMS record is represented as a tuple $\langle x, y, t, d \rangle$.

We construct a social graph $\mathcal{G} = \langle \mathcal{N}_T, \mathcal{E} \rangle$, based on the aggregated traffic of M months. We use \mathcal{N}_T to denote the set of mobile phone users that appear in the dataset. \mathcal{N}_T contains about 90 million unique cell phone numbers. Among the numbers that appear in \mathcal{N}_T , only some of them are clients of the mobile phone operator: we denote that set \mathcal{N}_O .

For this study, we had access to basic demographic information for a subset of the operator clients, that we denote \mathcal{N}_{GT} (where GT stands for *ground truth*). The size of this labeled set $|\mathcal{N}_{GT}|$ is about 500,000 users. The following relation holds between the three sets: $\mathcal{N}_{GT} \subset \mathcal{N}_O \subset \mathcal{N}_T$.

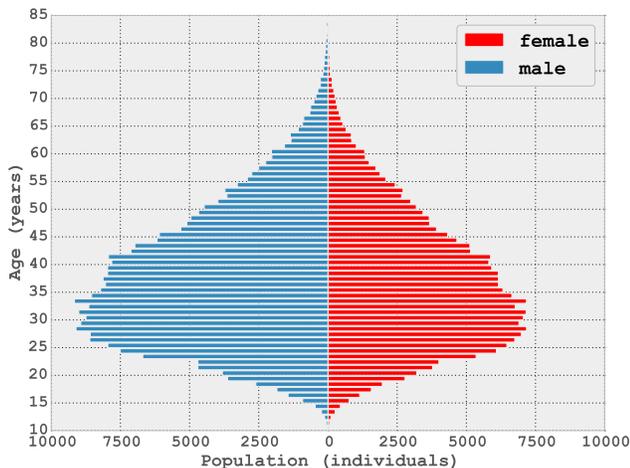


Fig. 1. Pyramid of ages of the labeled set \mathcal{N}_{GT} , showing the number of individuals for each age value (in years).

Fig. 1 shows the pyramid of ages of the labeled set. This pyramid is different from the age pyramid of the entire population, since it only contains mobile phone users (that pertain to \mathcal{N}_{GT}). Some basic observations on this pyramid: there are more men (56.83%) than women (43.17%) in the labeled set. The mean age is 37.23 years for men and 36.47 years for women.

III. OBSERVATIONAL STUDY

We report in this section the quantitative analyses that we performed on the dataset, which provide the basis for the prediction algorithms (which will be described in section IV). Our raw data input are the transaction logs (that contain billions of records at the scale of Mexico). A first step in the study was to generate characterization variables for each user, which summarize their individual and social behavior. We also describe the preprocessing performed on the data, the key features identified with PCA (Principal Component Analysis), and the differences observed. We will focus on some examples to illustrate the kind of observations we obtained for many variables. Statistically significant differences have been found, which motivates our attempt to identify gender and age based on communication patterns.

A. Characterization Variables

In this study we chose to characterize the users in \mathcal{N}_O (i.e. clients of the mobile operator), since for users in $\mathcal{N}_T \setminus \mathcal{N}_O$ we can only see part of their calls and messages (those exchanged with users in \mathcal{N}_O), thus including them would require a different calibration. For users in \mathcal{N}_O , we computed the following variables which characterize their calling consumption behavior (also called “behavioral variables” in [4]).

- *Number of Calls.* We consider incoming calls i.e., the total number of calls received by user u during a period of three months, as well as outgoing calls i.e., total number of calls made by user u . Additionally, we distinguish

whether those calls happened during the weekdays (Monday to Friday) or during the weekend; and we further split the weekdays in two parts: the “daylight” (from 7 a.m. to 7 p.m.) and the “night” (before 7 a.m. and after 7 p.m.). We thus have 3×4 variables for the number of calls, given by the Cartesian product: $[in, out, all] \times [weekdaylight, weeknight, weekend, total]$.

- *Duration of Calls.* We calculate the total duration of incoming calls and outgoing calls of user u during the period of three months. As before, we distinguish between weekdays (by daylight and by night) and weekends, to get a total of 12 variables for the duration of calls.
- *Number of SMS.* We consider incoming messages (received by user u) and outgoing messages (sent by user u). Similarly we distinguish between weekdays (by daylight and by night) and weekends, to get a total of 12 variables for the number of SMS.
- *Number of Contact Days.* We consider the number of days where the user has activity. We distinguish between calls and SMS, and between incoming, outgoing or any activity. This way we get 6 variables related to the number of activity days.

We also computed variables which characterize the social network of users based on their use of the cell phone (also called “social variables” in [4]).

- *In/Out-degree of the Social Network:* The in-degree for user u is the number of different phone numbers that called or sent an SMS to that user. The out-degree is the number of distinct phone numbers contacted by user u .
- *Degree of the Social Network:* The degree is the number of unique phone numbers that have either contacted or been contacted by user u (via voice or SMS).

B. Data Preprocessing

Many of the variables that we generated have a right skewed or heavy tailed distribution. Our experiments showed that this skewness affects the results given by Machine Learning algorithms (described in section IV-B). Therefore as part of the data preprocessing we also considered the logarithmic version of the variables. We discuss this preprocessing in more detail for one variable, that we use as running example: *in-time-total*, i.e. the total duration of incoming calls for a given user.

As can be seen in Table I (left), the quartiles of the variable *in-time-total* lie in different orders of magnitude, in particular the ratio $IQR/Q_2 = (Q_3 - Q_1)/Q_2$ is well above 1.

To improve the results given by the Machine Learning methods, we transform the data using the function $T(x) = \log_{10}(x + 1)$. After the transformation, we found the statistics in Table I (right). As we can see, the quartiles are in the same order of magnitude, and the ratio IQR/Q_2 is below 1. The resulting distribution is shown in Fig. 2.

In conclusion, we decided to include both plain variables as well as their logarithmic values, and let our Machine Learning algorithms select which variables are most relevant for modeling a given target variable (e.g. gender and age). We also rescaled all variables to take values between 0 and 1.

TABLE I
STATISTIC SUMMARY FOR *in-time-total* AND ITS LOGARITHMIC TRANSFORMATION

	<i>in-time-total</i> (seconds)	$\log(\textit{in-time-total} + 1)$
count	131770.00	131770.00
mean	16239.28	3.31
std	50023.16	1.23
min	0.00	0.00
25%	662.00	2.82
50%	3838.00	3.58
75%	14108.00	4.14
max	4045686.00	6.60

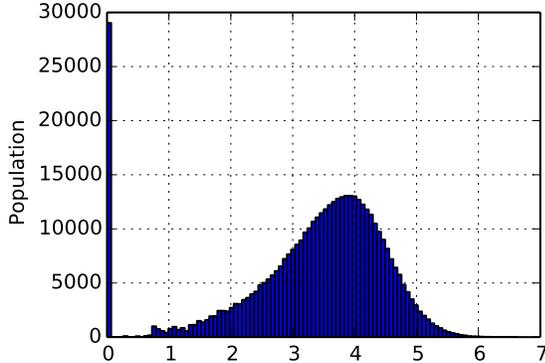


Fig. 2. Histogram of $\log(\textit{in-time-total} + 1)$. The distribution can easily be seen as a mixture between a distribution with all its density in 0 (users that have no incoming calls) and a bell shaped distribution for $\log(\textit{in-time-total} + 1) > 0$.

C. Insights on Key Features from PCA

We performed PCA (Principal Component Analysis) on the characterization variables, in order to gain information on which are the most important variables. This gave us interesting insights on the key features of the data. We describe the first 4 eigenvectors (which account for 89.6% of the variance).

The first eigenvector retains 76.0% of the total variance. This eigenvector is dominated by the logarithmic version¹ of the total number of calls, total duration of calls and total number of SMS. This result shows that the level of activity of users exhibits the highest variability, and therefore is a good candidate to characterize users' social behavior.

The second eigenvector, which retains 6.7% of the variance, gives high positive coefficients to “outgoing” variables (number of outgoing calls, duration of outgoing calls, number of SMS sent) and negative coefficients to “incoming” variables (number of incoming calls, duration of incoming calls, number of SMS received). This suggests that the difference of outgoing minus incoming communications is also a good variable to describe users' social behavior.

The third eigenvector (with 4.4% of the variance) gives positive coefficients to the “voice call” variables, and negative

¹We note that in all eigenvectors, the logarithmic version of the variables got systematically higher coefficients than the plain variables.

coefficients to the “SMS” variables (intuitively the difference between voice and SMS usage is relevant).

The fourth eigenvector (with 2.5% of the variance) gives positive coefficients to the “weeknight” and “weekend” variables (communications made non-working hours, i.e. during the night or during the weekend), and negative coefficients to “weeklight” variables (communications made during the day, from Monday to Friday, which correspond roughly to working hours).

D. Observed Gender Differences

We report in Table II the average of 6 key variables, distinguished by gender. We know from PCA that the number of calls, and the total duration of calls made by a user, characterize their level of activity. A natural question is whether differences can be observed between genders. We found that for all the variables, there is a significant difference between genders, with a very small p -value ($p < 10^{-10}$). Recall that these values are computed as the aggregation of calls during a period of $M = 3$ months for users in \mathcal{N}_{GT} (about 500,000 users).

TABLE II
SAMPLE MEAN OF KEY VARIABLES FOR FEMALE AND MALE USERS. THE DURATIONS ARE EXPRESSED IN SECONDS.

Variable	Female	Male
$\hat{\mu}(\text{total duration})$	10038.75	10663.17
$\hat{\mu}(\text{total duration outgoing})$	6359.96	7239.53
$\hat{\mu}(\text{total duration incoming})$	3678.78	3423.64
$\hat{\mu}(\text{number of calls total})$	72.847	81.348
$\hat{\mu}(\text{number of calls outgoing})$	44.136	50.047
$\hat{\mu}(\text{number of calls incoming})$	28.710	31.301

Table II shows that men have on average higher total number of calls, and total duration of calls (measured in seconds). However an interesting pattern can be seen when we distinguish incoming and outgoing calls: the duration of outgoing calls is higher for men, but the duration of incoming calls is higher for women. It follows that the net duration of calls (the difference between outgoing and incoming calls) has a marked gender difference: the sample mean $\hat{\mu}(\text{net total duration}) = 3815.88$ seconds for men, and $\hat{\mu}(\text{net total duration}) = 2681.18$ seconds for women. We note that the number of outgoing calls is higher than the number of incoming calls for both men and women, due to a particularity of our dataset (for all the users in \mathcal{N}_T the total number of incoming and outgoing calls is the same, but for users in \mathcal{N}_{GT} there is a higher proportion of outgoing calls).

We also compute the conditional probability $p(g'|g)$ that a random call made by an individual with gender g has a recipient with gender g' , where we denote male by M and female by F . For the calls originated by male users, we found that $p(F|M) = 0.3735$ and $p(M|M) = 0.6265$. For the calls originated by female users, we found that $p(F|F) = 0.4732$ and $p(M|F) = 0.5268$. We can see a difference between

genders, in particular men tend to talk more with men, and women tend to talk more with women. More precisely:

$$\begin{aligned} p(M|F) < p(M) = 0.5683 < p(M|M) \\ p(F|M) < p(F) = 0.4317 < p(F|F) \end{aligned} \quad (1)$$

Similar observations have been made in the case of the Facebook social graph [5].

E. Observed Age Differences

We approached the study of mobile phone usage patterns according to age by dividing the population in $C = 4$ categories: below 25 years, from 25 to 34 years, from 35 to 49 years, and 50 years or above. We use this same structure for age prediction (in section IV-C).

Since we are dealing with more than 2 groups, comparing differences between groups requires using the correct tool, as the probability of making a type I error (null hypothesis incorrectly rejected) increases. In order to compare the means (of the log) of the variables for each age group, we conduct a Tukey's HSD (Honest Significant Difference) test. This method tests all groups, pairwise, simultaneously. We found a list of 20 variables for which the null hypothesis of same mean (H_0) is rejected for all pair of groups, i.e. $\mu_i \neq \mu_j$ for every $i \neq j$.

TABLE III
TUKEY HSD FOR THE VARIABLE $\log_{10}(in-time-total + 1)$

group1	group2	meandiff	lower	upper	reject
0	1	0.1567	0.1328	0.1807	True
0	2	0.1326	0.1088	0.1564	True
0	3	0.2367	0.2122	0.2612	True
1	2	-0.0242	-0.0407	-0.0076	True
1	3	0.08	0.0625	0.0975	True
2	3	0.1041	0.0868	0.1214	True

We illustrate the difference between age groups for our running example: *in-time-total* (total duration of incoming calls per user). Table III show the result of Tukey HSD (where FWER=0.05) for the variable $\log_{10}(in-time-total+1)$, obtained after the preprocessing step. The 4 age groups are labeled 0, 1, 2, 3. Pairwise comparisons are done for all combinations (of group1 and group2). The null hypothesis is rejected for all pairs; in other words, all the groups are found statistically different respect to this variable.

In Fig. 3 we plot the distribution of $\log_{10}(in-time-total + 1)$ for different pairs of age groups. The following results can be observed from the plots:

- The distribution for the group of people aged *over 50* is shifted to the right in comparison with all the other age groups. This implies that people from this age group do talk more when called than people from any other age group. Figures 3c, 3e, 3f.
- The distribution for the group of people aged *below 25* is shifted to the left. This distribution shows less kurtosis and a higher variance, meaning that this population is

more spread in different levels of $\log_{10}(in-time-total+1)$. Figures 3a, 3b, 3c.

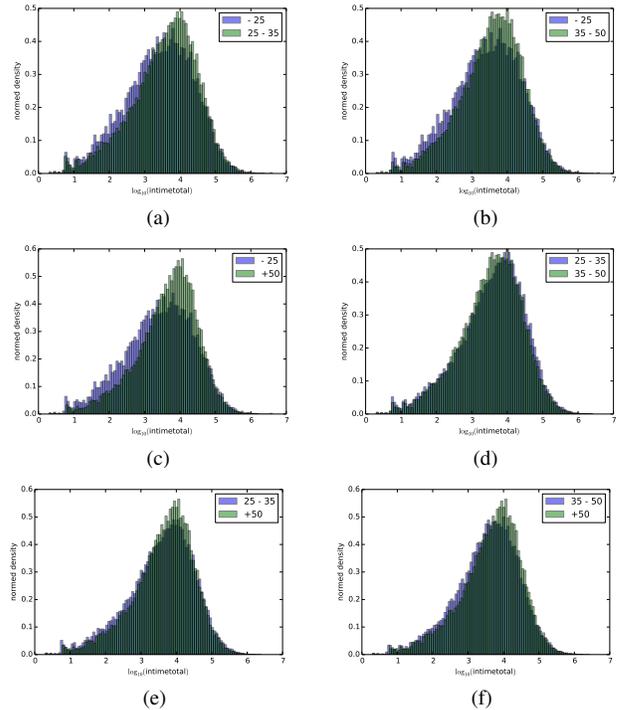


Fig. 3. Pairwise comparison of $\log_{10}(in-time-total+1)$ between different age groups. The plots show the normed density for $\log_{10}(in-time-total+1) > 0$.

F. Links between Age Groups

We study here the links between users, according to their age. Fig. 4 shows the matrix $C_{i,j}$ that contains the number of links between users of age i and age j . For each user $u \in \mathcal{N}_{GT}$, we compute the number of direct contacts of u that belong to \mathcal{N}_{GT} and have age j . We sum over all the users of age i to get the number $C_{i,j}$. As we can see in the figure, the diagonal of the matrix has clearly higher values than the rest of it, meaning that users are more likely to establish communications with someone of their own age. This strong *age homophily* has also been observed in [5], and in smaller social networks [6].

The communication preferences can also be seen in Fig. 5, which shows the number of links according to the age difference between users. The highest number of links is observed when the difference is $\delta = 0$. The number of links decreases with the age difference, except around the value $\delta = 21$, where an interesting inflection point can be observed; possibly relating to different generations (e.g. parents and children).

IV. AGE AND GENDER PREDICTION

This section describes the models that we used to estimate the age and gender of users found in the dataset $\mathcal{N}_O \setminus \mathcal{N}_{GT}$. We show the results obtained using standard Machine Learning models based on node attributes, applied to the prediction of gender (section IV-B) and age (section IV-C). We introduce a novel algorithm that leverages the communication network

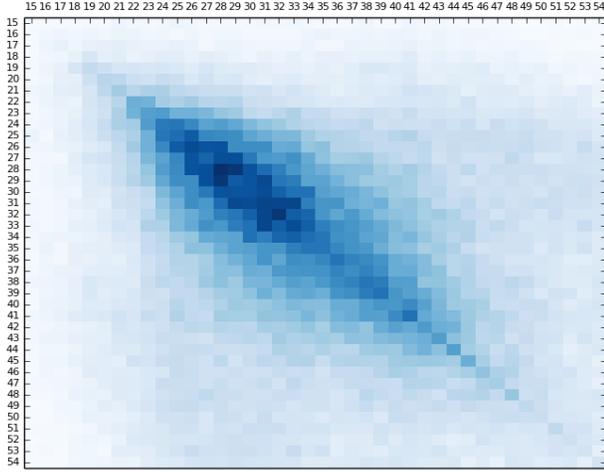


Fig. 4. The matrix $C_{i,j}$ of communications between users of age i and age j (the ages are expressed in years).

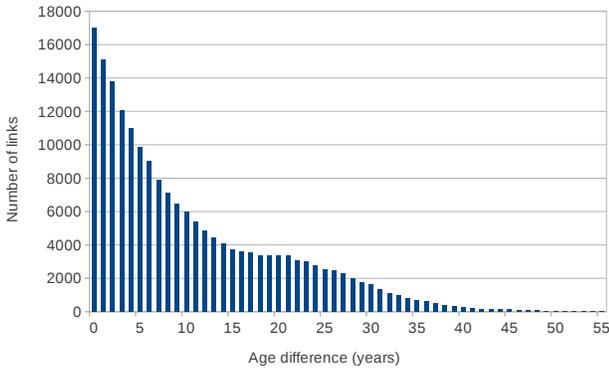


Fig. 5. Number of links as a function of the age difference between users.

topology to generate age predictions (section IV-D), and finally show how it can be combined with the Machine Learning models (section IV-E).

We note that the feature variables are known for the whole set \mathcal{N}_O , while the target variables age and gender are known only for users in the set \mathcal{N}_{GT} . We therefore use nodes belonging to \mathcal{N}_{GT} for our training and validation set, to predict both age and gender for the remaining users in \mathcal{N}_O .

A. Population Pyramid Scaling (PPS)

In the following subsections, we will present algorithms for gender and age prediction which generate for each node, a probability vector over the possible categories (gender or age groups). At the end of their execution, when we “observe” the system, we are required to collapse the probability vectors to a specific gender or age group. Choosing how to perform this collapse is not an obvious matter. In effect, we want to collapse the probability state of the system as a whole and not for each node independently. In particular, we want to impose external constraints on our solution, namely that the gender or age group distribution for the whole network be that of the ground truth. To achieve this, we developed a method that we

Algorithm 1: Population Pyramid Scaling

```

1 foreach node  $i$  and group  $k$  do
2   compute the probability  $p_{i,k}$  that node  $i$  belongs to
   group  $k$  using the unconstrained algorithm;
3 end
4 Create an ordered list  $T$  of tuples  $(i, k, p_{i,k})$ ;
5 Sort list  $T$  in descending order by the column  $p_{i,k}$ . The
   list  $T$  will be iterated starting with the element with the
   highest probability;
6 foreach element  $(i, k, p_{i,k}) \in T$  do
7   if node  $i$  has not been assigned to a group then
8     if less than  $N_k$  nodes assigned to group  $k$  then
9       assign node  $i$  to group  $k$ ;
10    end
11  end
12 end

```

call *Population Pyramid Scaling*. This algorithm takes, as a hyper-parameter, the proportion q of nodes to be predicted. For example, we use $q = 1/2$ to generate predictions for the 50% of nodes which got better classification results from the unconstrained method.

The PPS procedure is described in Algorithm 1. Note that the population to predict has size $N = q \times |\mathcal{N}_O|$. For each category k we compute the number of nodes N_k that should be allocated to category k in order to satisfy the distribution constraint (the gender or age distribution of \mathcal{N}_{GT}), and such that $\sum_{k=1}^C N_k = N$ (where C is the number of categories or groups).

B. Gender Prediction

For gender prediction, several algorithms were evaluated, with a preference for algorithms more restrictive respect to the functions that they adjust. Some of the algorithms used are: Naive Bayes, Logistic Regression, Linear SVM, Linear Discriminant Analysis and Quadratic Discriminant Analysis. As previously described, as part of data preprocessing, log transformation of the variables are added and the values are standardized to the $[0, 1]$ interval.

TABLE IV
BEST CLASSIFIERS CONFIGURATION FOR GENDER PREDICTION

Algorithm	Best configuration
LinearSVC	dual = True; penalty = L2; loss = L1; C = 1; k = 100; training set = 200,000
LogisticRegression	penalty = L1; C = 10; k = 100; training set = 200,000

The best results were obtained with Linear SVM and Logistic Regression. Table IV summarizes the classifiers configuration. To find these parameters, we used grid search over a predefined set of parameters. For instance, the parameter C for Logistic Regression takes its values in the

set $\{0.1, 0.3, 1, 3, 10\}$. Different number of attributes were evaluated before training the model ($k \in \{10, 30, 100\}$). The labeled nodes were split in a training set (70%) and a validation set (30%).

TABLE V
PRECISION OBTAINED FOR GENDER PREDICTION

Parameter q	1	1/2	1/4	1/8
Accuracy	66.3%	72.9%	77.1%	81.4%

After performing PPS (to ensure the correct proportion of men and women), we obtained the results shown in Table V. As expected, the accuracy of our predictions improve when we decrease the parameter q , which provides a trade-off between precision and coverage. We reach a precision of 81.4% when tagging 12.5% of the users.

In the following paragraphs, we briefly recall some details of the classifiers that gave the best results, in order to clarify the meaning of the configuration parameters in Table IV. Those parameters are the ones required by the Scikit-learn library [7]. In addition, Pandas [8] and Statsmodels [9] have been used for the exploratory analysis.

1) *Linear SVC*: Classification using Support Vector Machines (SVC) requires optimizing the following function [10]:

$$\min_{\omega} \frac{1}{2} \omega^T \omega + C \sum_{i=1}^l \xi(\omega, x_i, y_i) \quad (2)$$

where y_i is the gender value and x_i is the vector of normalized observed variables; ω describes the hypothesis function, and C is the regularization parameter. In our case, we choose to use $\xi(\omega, x_i, y_i) = \max(1 - y_i, \omega^T x_i, 0)$. This setup is called L1-SVM as $\xi(\cdot)$ defines the loss function. Implementation details can be found in [7], [10], [11].

2) *Logistic Regression*: A standard way to estimate discrete choice models is using index function models. We can specify the model as $y^* = \omega x + \epsilon$, where y^* is an unobserved variable. We use the following criteria to make a choice:

$$\begin{cases} y = 1 = \text{Female} & \text{if } y^* > 0 \\ y = 0 = \text{Male} & \text{if } y^* < 0 \end{cases}$$

Additionally, we use L1 regularization (for feature selection and to reduce overfitting). The complete formulation is to optimize:

$$\min_{\omega} \|\omega\|_1 + C \sum_{i=1}^l \xi(\omega, x_i, y_i)$$

where $\xi(\omega, x_i, y_i) = \log(1 + \exp(-y_i \omega^T x_i))$. For more information refer to [7], [11], [12].

C. Age Prediction using Node Attributes

We now tackle the problem of detecting the age of the users using the properties of the nodes (users). As a first approach, we use the Machine Learning armory to perform the detection. In order to reduce the complexity of the target variable, we

partition it into C age categories ($C = 4$): below 25 years old, from 25 to 34 years old, from 35 to 49 years old, and 50 years old or above (as in section III-E). Given this set of categories, we found the best results using Multinomial Logistic (MNLogistic). This method is a generalization of Logistic Regression for the case of multiples categories (refer to [9], [12]).

A problem we encountered when using MNLogistic is the overfitting to the categories with higher frequencies, in this case the classification in categories 25 to 34 years and 35 to 49 years of more elements than expected. In effect, in our Training Set, the age groups have the following distribution:

Age group	10-25	25-35	35-50	50+
Population	12.1%	35.45%	37.45%	15%

But when using the MNLogistic algorithm, we obtained predictions with the following distribution:

Age group	10-25	25-35	35-50	50+
Population	0.66%	52.97%	45.52%	0.84%

To solve this issue we used the PPS (Population Pyramid Scaling) method of section IV-A. After performing PPS, we obtained the results presented in Table VI (in section IV-F).

D. Age Prediction using Network Topology

As discussed in section III-F, there is a strong age homophily among nodes in the communication network, yet the Machine Learning algorithms we have employed so far are mostly blind to this information, that is, their predictive power relies solely on user attributes ignoring the complex interactions given by the mobile network they participate in.

In this section we propose an algorithm which can harness the information given by the structure of the mobile network and in this way, leverage hidden information as the homophily patterns (shown in Fig. 4).

1) *Communication Network Structure*: We briefly describe how we construct the communication network \mathcal{G} . To each user (phone number) we assign a node in the network, and to each pair of users x and y communicating via voice calls or SMS we assign a link (noted $x \sim y$). We can also assign a weight $w_{x,y}$ to each link between user x and user y , expressing the number of communications, the total activity or any other property of the link. In the present study we choose to set $w_{x,y} = 1$ if there is any kind of communication between nodes x and y , and $w_{x,y} = 0$ otherwise. We do this both to reduce the complexity of the analysis, and because as a first approach, we are interested in exploring how the bare network topology enhances our predictions of the target variable.

2) *Reaction-Diffusion Algorithm*: In our dataset, we have the values for the target variable (age) of the nodes in \mathcal{N}_{GT} , and we can also see from Fig. 4 that neighbouring nodes are more likely to belong to the same age category. A mathematical model that can take advantage of this information together with the topology of the network to infer the target

values for the remaining nodes (in $\mathcal{N}_O \setminus \mathcal{N}_{GT}$) is that of a diffusion process in a graph. At each time step the information (value of target variable at a given node) is diffused or spread to its neighbours. In this way, given enough time steps, the information from nodes in \mathcal{N}_{GT} is diffused to the entire network. Now, if $|\mathcal{N}_{GT}| \ll |\mathcal{N}_O|$, which is the case in our study, pure diffusion may not be strong enough to have the information in \mathcal{N}_{GT} significantly affect the values of the target variables in the entire network. To remedy this, we include a reactive term in our algorithm where at each time step, nodes in \mathcal{N}_{GT} are reinforced with their value at time $t = 0$. Given that we have partitioned our target variable in C categories, we also found it advantageous to model our reaction-diffusion process as one where the information being diffused is the probability distribution for each node (noted $g_{x,t}$) to belong to each category. We detail the algorithm below.

For each user x we define the initial state $f_x \in \mathbb{R}^C$ (where C is the number of categories) as having components

$$(f_x)_i = \begin{cases} \delta_{i,a(x)} & \text{if } x \in \mathcal{N}_{GT} \\ 1/C & \text{if } x \notin \mathcal{N}_{GT} \end{cases} \quad (3)$$

where $a(x)$ is the age category of user x and $\delta_{i,a(x)}$ is the Kronecker delta function. Then we define $g_{x,t}$ as

$$g_{x,0} = f_x$$

$$g_{x,t} = (1 - \lambda) f_x + \lambda \frac{\sum_{x \sim y} w_{y,x} g_{y,t-1}}{\sum_{x \sim y} w_{y,x}} \quad (4)$$

where $x \sim y$ means that there is a link between x and y ; $w_{y,x}$ is the weight of the link; and λ is a hyper-parameter which tunes the relative strength of the reinforcement and diffusion terms (which we set to $\lambda = 0.5$ in our experiments). Note that $g_{x,t} \in \mathbb{R}^C$ is the discrete probability measure for node x at time t , in particular the sum $\sum_{i=1}^C (g_{x,t})_i = 1 \forall x, t$.

As previously stated, these equations are similar to those of a *reaction-diffusion* process (we note that it can also be seen as a Jacobi method for the appropriate linear system). For the experimental results, we consider a simpler model by taking $w_{y,x} = 1 \forall x \sim y$ and 0 otherwise. The equation for $g_{x,t}$ becomes

$$g_{x,t} = (1 - \lambda) f_x + \lambda \frac{\sum_{x \sim y} g_{y,t-1}}{|\{y : x \sim y\}|}. \quad (5)$$

We iterate this process m times (for $1 \leq t \leq m$). In our experiments, $m = 30$ was sufficient for the process to converge.

For each x , we obtain a vector $g_x = g_{x,m}$. With this we can get a prediction for the age of x given by $\operatorname{argmax}_{1 \leq i \leq 4} (g_x)_i$. This prediction, in contrast with the prediction performed the MNLogistic model (based on node attributes), gave us a population pyramid closer to the ground truth:

Age group	10-25	25-35	35-50	50+
Population	7.26%	32.42%	50.49%	10.36%

After adjusting the distribution with the PPS algorithm (from section IV-A), we obtained the results shown in Table VI.

E. Enriching the Graph Algorithm with Node Attributes

We propose here an algorithm to predict the age of users that leverages the PPS algorithm, the node classification of section IV-C and the pure graph-based *Reaction-Diffusion* algorithm. We define as initial state:

$$f_x = \begin{cases} \delta_{i,a(x)} & \text{if } x \in \mathcal{N}_{GT} \\ \text{ML}(x) & \text{if } x \notin \mathcal{N}_{GT} \end{cases} \quad (6)$$

where $\text{ML}(x)$ is the result given by the best *Machine Learning* algorithm of section IV-C (i.e. Multinomial Logistic).

Then, as before, the iterative process follows Equation (5). In this case, the hyper-parameter λ provides a trade-off between the information from the network topology and the initial information obtained with Machine Learning methods over node attributes (here again we take $\lambda = 0.5$).

F. Summary of Results

Finally, Table VI summarizes the results obtained with the different methods: Machine Learning (ML) alone, Reaction-Diffusion (RDif) alone, and the combined method (ML + RDif). We report for each case the accuracy obtained, that is the percentage of correct predictions on the validation set.

TABLE VI
PRECISION OBTAINED FOR AGE PREDICTION

Population	ML	RDif	ML + RDif
$q=1$	36.9%	43.4%	38.1%
$q=1/2$	42.9%	47.2%	46.3%
$q=1/4$	48.4%	56.1%	52.3%
$q=1/8$	52.7%	62.3%	57.2%

The table shows that taking a smaller q improves the accuracy of the results. Our experiments also show that the RDif (Reaction-Diffusion) algorithm outperforms the ML predictions based on node attributes. It is also interesting to remark that the RDif algorithm outperforms the combined method. The best precision obtained is 62.3% of correctly predicted nodes, when tagging 12.5% of the population. Note that random guessing the age group (between 4 categories) would yield a precision of 25%.

V. CONCLUSION AND FUTURE WORK

To our knowledge, this work provides the first extensive study of social interactions in the country of Mexico focusing on gender and age, based on mobile phone usage. From a sociological perspective, the ability to analyze the communications between tens of millions of people allows us to make strong inferences and detect subtle properties of the social network.

As described in section III, the graph we constructed has very rich link semantics, containing a detailed description of the communication patterns (45 characterization variables). With PCA, we found that most of the variance of the characterization variables is contained in a low dimensional subspace. Motivated by these results, we focused on how the statistical properties of the most informative attributes vary with both

gender and age. In section III-D, we make two interesting observations: (i) there is a gender homophily in the communication network (see Equation 1); (ii) an asymmetry respect to incoming and outgoing calls can be observed between men and women, possibly reflecting a difference of roles in Mexican society (it would be interesting to see how these differences change in other regions like Europe or the United States). We also compared communication habits for different age groups, and found statistically significant differences. Finally, our most important observational contribution is the study of correlations between age groups in the communication network, as summarized in Fig. 4 and 5. We observe a strong age homophily [6], and a strong concentration of communications centered around the age interval between 25 and 45 years. But we also notice weaker modes in both figures, which raise interesting sociological questions (e.g. whether they reflect a generational gap).

The second key contribution of this work was to study and propose novel methods to infer the gender and age of users in the mobile network. As a first approach, described in sections IV-B and IV-C, we used a set of standard Machine Learning tools finding that Logistic Regression and Linear Support Vector Machine algorithms gave us the best results. However, these techniques cannot harness the topological information of the network to explore possible correlations between the users' age groups. To leverage this information, we proposed an purely graph based algorithm inspired in a *reaction-diffusion* process, and demonstrated that with this methodology we could predict the age category for a significant set of nodes in the network. Our experiments showed that the *reaction-diffusion* method provides the best predictive power on a real-world large scale dataset.

There are multiple directions in which this work can be extended. We highlight the following:

a) *Analysis of Hyper Parameters*: The analysis of the prediction performance as a function of the hyper-parameters q and λ , used in sections IV-C and IV-D, is important for a fine tuning of the algorithm. We are also interested in studying the effect of variations in the weights $w_{x,y}$ used in the diffusion process (e.g. use the intensity of communication or the geolocation data to weight the links). In particular, we want to explore how the network topology information can be combined with nodes features to improve the joined (ML + RDif) methodology.

b) *Extend Depth*: A statistics quasi-experiment can be built from this method [13]. In this case, we want to know whether the differences in the observed behavior can be accounted to gender and age, or are consequences of differences in the ego-network induced by phone calls. This quasi-experiment can be performed using Propensity Score [14], and may provide sociological insights.

c) *Extend Width*: One direction that we are currently investigating is to apply the methodology presented here to predict variables related to the users' spending behavior. In [15] the authors show correlations between social features and spending characterizations, for a small population (52

individuals). We are interested in applying our methodology to predict spending behavior characteristics on a much larger scale (millions of users).

Another research direction is to use the geolocation information contained in the Call Details Records. Recent studies have focused on the mobility patterns related to cultural events –for instance sport related events [16], [17]– which might exhibit differences between genders and age groups. Looking at mobility patterns through the lens of gender and age characterization will provide new features to feed the Machine Learning part of our methodology, and more generally will provide new insights on the human dynamics of different segments of the population.

REFERENCES

- [1] J. Blumenstock and N. Eagle, "Mobile divides: gender, socioeconomic status, and mobile phone use in Rwanda," in *Proceedings of the 4th ACM/IEEE International Conference on Information and Communication Technologies and Development*. ACM, 2010, p. 6.
- [2] J. E. Blumenstock, D. Gillick, and N. Eagle, "Who's calling? demographics of mobile phone use in Rwanda," *Transportation*, vol. 32, pp. 2–5, 2010.
- [3] E. G. Katz and M. C. Correia, *The economics of gender in Mexico: Work, family, state, and market*. World Bank Publications, 2001.
- [4] V. Frias-Martinez, E. Frias-Martinez, and N. Oliver, "A gender-centric analysis of calling behavior in a developing economy using call detail records," in *AAAI Spring Symposium: Artificial Intelligence for Development*, 2010.
- [5] J. Ugander, B. Karrer, L. Backstrom, and C. Marlow, "The anatomy of the Facebook social graph," *structure*, vol. 5, p. 6, 2011.
- [6] M. McPherson, L. Smith-Lovin, and J. M. Cook, "Birds of a feather: Homophily in social networks," *Annual review of sociology*, pp. 415–444, 2001.
- [7] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, and V. Dubourg, "Scikit-learn: Machine learning in python," *The Journal of Machine Learning Research*, vol. 12, p. 2825–2830, 2011.
- [8] W. McKinney, "Data structures for statistical computing in python," in *Proceedings of the 9th Python in Science Conference*, S. van der Walt and J. Millman, Eds., 2010, pp. 51 – 56.
- [9] J. Seabold and J. Perktold, "Statsmodels: Econometric and statistical modeling with python," in *Proceedings of the 9th Python in Science Conference*, 2010.
- [10] C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and S. Sundararajan, "A dual coordinate descent method for large-scale linear SVM," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, p. 408–415.
- [11] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: a library for large linear classification," *The Journal of Machine Learning Research*, vol. 9, p. 1871–1874, 2008.
- [12] W. H. Greene, *Econometric Analysis*, 7th ed. Prentice Hall, Feb. 2011.
- [13] W. Shadish, T. D. Cook, and D. T. Campbell, *Experimental and quasi-experimental designs for generalized causal inference*. Wadsworth Cengage learning, 2002.
- [14] P. Rosenbaum and D. Rubin, "The central role of the propensity score in observational studies for causal effects," *Biometrika*, vol. 70, no. 1, pp. 41–55, 1983.
- [15] V. K. Singh, L. Freeman, B. Lepri, and A. S. Pentland, "Predicting spending behavior using socio-mobile features," in *Social Computing (SocialCom), 2013 International Conference on*. IEEE, 2013, pp. 174–179.
- [16] N. Ponieman, A. Salles, and C. Sarraute, "Human mobility and predictability enriched by social phenomena information," in *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. ACM, 2013, pp. 1331–1336.
- [17] F. H. Xavier, C. H. Malab, L. Silveira, A. Ziviani, J. Almeida, and H. Marques-Neto, "Understanding human mobility due to large-scale events," in *Third International Conference on the Analysis of Mobile Phone Datasets (NetMob)*, 2013.