



# On the Unicity of Smartphone Applications

Jagdish Prasad Achara, Gergely Acs, Claude Castelluccia

## ► To cite this version:

Jagdish Prasad Achara, Gergely Acs, Claude Castelluccia. On the Unicity of Smartphone Applications. 2015. hal-01181040v1

**HAL Id: hal-01181040**

**<https://inria.hal.science/hal-01181040v1>**

Preprint submitted on 29 Jul 2015 (v1), last revised 29 Oct 2015 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On the Unicity of Smartphone Applications

Jagdish Prasad Achara  
INRIA  
jagdish.achara@inria.fr

Gergely Acs  
INRIA  
gergely.acs@inria.fr

Claude Castelluccia  
INRIA  
claude.castelluccia@inria.fr

## ABSTRACT

Prior works have shown that the list of apps installed by a user reveal a lot about user interests and behavior. These works rely on the semantics of the installed apps and show that various user traits could be learnt automatically using off-the-shelf machine-learning techniques. In this work, we focus on the re-identifiability issue and thoroughly study the unicity of smartphone apps on a dataset containing 54,893 Android users collected over a period of 7 months. Our study finds that any 4 apps installed by a user are enough (more than 95% times) for the re-identification of the user in our dataset. As the complete list of installed apps is unique for 99% of the users in our dataset, it can be easily used to track/profile the users by a service such as Twitter that has access to the whole list of installed apps of users. As our analyzed dataset is small as compared to the total population of Android users, we also study how unicity would vary with larger datasets. This work emphasizes the need of better privacy guards against collection, use and release of the list of installed apps.

## 1. INTRODUCTION

People are all unique the way they are or they look. They can be easily identified from their DNA sequences, fingerprints, Iris scans, web browsers and so on. Also, a combination of various attributes, such as their age, address or religion [14] makes them unique. Recently, some studies have shown that people are also unique in the way they behave. For example, de Montjoye *et al.* [2] illustrated this *behavioural uniqueness* by showing that people are unique in the way they move. In fact, they show that only four spatio-temporal positions are enough to uniquely identify a user 95% of the times in a dataset of one and a half million users [2]. It has also been shown recently that people are unique in the way they consume and purchase goods online [3]. Similarly, other studies showed that people are unique in the way they configure their browser or browse the web [4] [11].

As smartphones have been widely deployed today all over the world and the list of installed/running applications (apps) on them is readily available to be accessed, the threat in terms of user privacy is huge if this data is collected, used and released without sufficient diligence in terms of privacy. This threat comes in two flavors: first, the semantics of the installed apps can tell a lot about the users' habits and interests [13], and second, the unicity of installed apps could make a user re-identifiable if this dataset is released. In fact, regarding the first privacy threat, [13] showed that user traits such as religion, relationship status, spoken languages, countries of interest, and whether or not the user is a parent of small children, can be easily predicted from the list or even the categories of the installed apps on smartphones. In this paper, we focus on the second privacy threat and measure the unicity of installed apps to be able to measure the risk of re-identification if app datasets are released in public or shared between two entities.

It is quite in the news these days<sup>1</sup> that Twitter has started to collect the list of apps that a user has installed. They claim to use this information for targeted interest-based advertising among others. However, it might be a privacy concern if Twitter shares this list of installed apps with an advertising company, even in pseudo-anonymized form, i.e., after removing all direct user identifiers (and even if app names are replaced with their hashes). This is because the advertising company might implicitly know a subset, say  $K$ , of installed apps of a user in which their ad library is present. So if  $K$  apps are enough to uniquely identify a user in the dataset, the advertiser would be able to re-identify the user in the Twitter dataset, and hence, learn about all the other installed apps of that user. By knowing this whole list of installed apps of a user, the advertiser can learn about that user's interests and habits (as demonstrated in [13]), and consequently, might be able to deliver the targeted ads directly in these apps in which its library is present. We believe that this is a real privacy threat to smartphone users (both Android and iOS) today as apps running on these OSs can access the list of installed/running apps. It is to be noted that Android apps do not require any permission to access the whole list of installed apps. On iOS, Apple does not provide a public API to access the list of installed apps but apps can get the list of currently running apps at any time. And if an app makes a frequent scan of currently running

<sup>1</sup><http://recode.net/2015/06/10/twitter-advertisers-can-now-target-you-based-on-the-other-apps-on-your-phone/>

apps over a period of time, the list can converge very fast to the list of installed apps.

**Contributions:** The contributions of our work are as follows.

- We show that 99% of the lists of installed application of users are unique out of a total of 55 thousands users. Moreover, as few as two applications are sufficient for an adversary to identify an individual's application list with a probability of 0.75 in our dataset. The re-identification probability increases to almost 0.95 if the adversary knows 4 apps. We stress that these results were obtained without considering any system apps (which are common for all users), and apps were identified only by the hash of their names. Incorporating additional information into their identifiers, such as app version, time of installation, etc., would increase these probabilities even more.
- We propose an unbiased estimate of the real uniqueness of any subset of applications, i.e., the probability that a randomly selected subset of apps with cardinality  $K$  is unique in the dataset. For this purpose, we use a Markov Chain Monte Carlo method to sample subsets of applications from a dataset uniformly at random. We prove that this chain is generally fast-mixing with most practical datasets, i.e., has a running time complexity which is roughly linear in the dataset size and  $K$ . This result might be of independent interest, as this technique can be used to sample subsets with arbitrary cardinality from any set-valued dataset.
- We attempt to predict the uniqueness of lists of applications in larger datasets using standard non-linear regression models. Although our model performs well on our limited app dataset as well as on mobility data with sufficiently large number of users, we conclude that our app dataset at hand is probably too small to accurately predict the uniqueness in a larger dataset such as the population of all Android users worldwide.

## 2. UNICITY AS A MEASURE OF RE-IDENTIFIABILITY

Let  $\mathbb{A}$  denote the universe of all apps, where each application is represented by a unique identifier in  $\mathbb{A}$ . A dataset  $D \subseteq 2^{\mathbb{A}} \setminus \{\emptyset\}$  is the ensemble of all apps of some set of individuals, where  $|D|$  denotes the number of individuals in  $D$ . A record  $D_u$ , which is a non-empty subset of  $\mathbb{A}$ , refers to all apps of an individual  $u$  in  $D$ . A set of applications with cardinality  $K$  is shortly called  $K$ -apps henceforth. The set of all  $K$ -apps over  $\mathbb{A}$  is denoted as  $\mathbb{A}^K$ .

**Definition 1 (Unicity)** Let  $\text{supp}(x, D)$  denote the support of  $x \in \mathbb{A}^K$  in  $D$ , i.e., the number of records in  $D$  which contain  $x$ . Then,

$$H_1 = \frac{|\{x : x \in \mathbb{A}^K \wedge \text{supp}(x, D) = 1\}|}{|\{x : x \in \mathbb{A}^K \wedge \text{supp}(x, D) \geq 1\}|}$$

is defined as the unicity (or uniqueness) of  $K$ -apps in dataset  $D$ .

The unicity of  $K$ -apps is the relative frequency of  $K$ -apps which are contained by only a single record. In general, *relative abundance distribution* (RAD)<sup>2</sup> is a relative frequency histogram  $\mathbf{H} = (H_1, H_2, \dots, H_n)$  of  $K$ -apps with respect to a dataset  $D$ , where  $H_i$  denotes the relative frequency of  $K$ -apps which are contained by exactly  $i$  records in  $D$ , i.e.,

$$H_i = \frac{|\{x : x \in \mathbb{A}^K \wedge \text{supp}(x, D) = i\}|}{|\{x : x \in \mathbb{A}^K \wedge \text{supp}(x, D) \geq 1\}|}.$$

Unicity is strongly related to re-identifiability, and we use it as a measure of privacy in this paper: it is the probability that the adversary, who only knows  $K$  applications installed on a user's device, can single out the record of this user in  $D$ . Indeed, any  $K$ -apps which is unique in  $D$  can be used as a personal identifiable information (PII) of its record owner. Specifically, if the adversary knows such  $K$ -apps, it can easily identify the corresponding record and retrieve all the applications installed by its owner, even if  $D$  is pseudo-anonymized (i.e., does not contain any direct PII such as device ID or personal name). Therefore, large unicity usually indicates a serious privacy risk in practice.

## 3. APPROXIMATING UNICITY WITH SAMPLING

To compute unicity, and RAD in general, the support of all different  $K$ -apps in  $D$  should be calculated. This is usually prohibitively expensive in practice. Rather, similarly to previous works [3, 2], we rely on sampling to estimate unicity. In particular, let  $\Omega^K$  denote the set of all  $K$ -apps which occur in at least one individual's record, i.e.,  $\Omega^K = \{x : x \in \mathbb{A}^K \wedge \text{supp}(x, D) \geq 1\}$ . We randomly sample a set  $V$  of  $K$ -apps from  $\Omega^K$ , and approximate the real unicity  $H_1$  by the sample unicity  $\hat{H}_1 = \frac{|\{x : x \in V \wedge \text{supp}(x, D) = 1\}|}{|V|}$ , where  $V \subseteq \Omega^K$  is the sample set, and  $n = |V|$  is the sample size.

### 3.1 Biased vs. unbiased estimation of unicity

How should we sample  $K$ -apps from the dataset? A popular technique, which has been used in several works [3, 2], first samples a user uniformly at random in  $D$ , and then a set of  $K$  applications from this user's record also uniformly at random. However, this simple technique provides a *biased estimation* of the unicity in Definition 1, if the estimator remains the sample mean  $\hat{H}_1$ , since  $E[\hat{H}_1] \neq H_1$ . Indeed, any  $K$ -apps which occur in more records than others also become more likely to be selected by this approach (assuming records have similar sizes). As a result, this sampling method is biased towards more popular  $K$ -apps, and the measured unicity is an underestimation of the real unicity  $H_1$  what one would get with an unbiased estimator of  $H_1$ . Such an unbiased estimator can be the sample unicity  $\hat{H}_1$  of  $K$ -apps which are sampled truly uniformly at random from  $D$ . This is also illustrated in Figure 4, where the sample unicity of biased and unbiased (i.e., uniform) samples are reported.

Before describing our unbiased estimation of unicity  $H_1$ , we shed some light on the privacy semantics behind the two sampling approaches. The biased technique approximates

<sup>2</sup>This term is often used in the field of ecology to describe the relationship between the number of observed species as a function of their observed abundance.

the success probability of an adversary who is more likely to know popular  $K$ -apps from the application set of any user. For instance, continuing the case of the advertiser from Section 1, the advertiser’s library is more likely to be used by popular apps (such as Facebook, Twitter, etc.), which are installed on many devices, rather than by other less popular apps. In general, the adversary can always learn certain  $K$ -apps of some users (e.g., the adversarial library is favoured by certain applications) with larger probability, which changes the sampled unicity  $\hat{H}_1$  accordingly; more popular apps tend to increase  $\hat{H}_1$  while unpopular apps tend to decrease it yielding a biased estimation of  $H_1$ . In the rest of the paper, we assume that the adversary can learn *any*  $K$ -apps of *any* users in  $D$  with equal probability, which is the most general assumption in practice. Therefore, we are interested in an unbiased estimator of  $H_1$ .

### 3.2 Uniform sampling of $K$ -apps

A unbiased estimation of  $H_1$  is obtained, if  $\hat{H}_1$  is computed over a sample set where each  $K$ -apps can appear with equal probability. Hence, our task is to sample an element from  $\Omega^K$  uniformly at random for any  $K$ . A first (naive) approach could be to use rejection sampling, i.e., sample a candidate  $K$ -apps from  $\mathbb{A}^K$  uniformly at random, and then accept this candidate as a valid sample only if it also occurs in  $D$ . Otherwise, repeat the process until a candidate is accepted. Although sampling a candidate from  $\mathbb{A}^K$  is straightforward, it is very likely to be non-existent in  $D$  (especially if  $K$  is large), and hence, its running complexity is  $O(|\mathbb{A}|^K)$  in the worst case. An alternative approach could be to enumerate  $\Omega^K$ , and choosing one element directly from  $\Omega^K$  uniformly at random. However, the complexity of this approach is still  $O(|D|(\max_u |D_u|)^K/K!)$ . Unfortunately, these naive methods provide acceptable performance only if  $K$  is small. As Table 1 shows, in our dataset,  $|\mathbb{A}| = 92210$ ,  $\max_u |D_u| = 541$ ,  $|D| = 54893$ , and we wish to estimate the unicity when  $1 \leq K \leq 10$ .

We instead propose a sampling technique based on the Metropolis-Hastings algorithm [10, 1], which is a Markov Chain Monte Carlo (MCMC) method. Our proposal has a worst-case complexity of only  $O(K|D|/H_1^*)$ , where  $H_1^*$  is roughly the unicity of  $K$ -apps in  $D$ . As the unicity of  $K$ -apps is large, especially if  $K$  is large, the complexity is approximately  $O(K|D|)$  in practice. Hence, our sampling technique remains reasonably fast even for larger values of  $K$ .

In particular, we construct an ergodic Markov chain, denoted by  $\mathcal{M}$ , such that its stationary distribution  $\pi$  is exactly the distribution that we want to sample from, that is, the uniform distribution over  $\Omega^K$ . Each  $K$ -apps in  $\Omega^K$  corresponds to a state of  $\mathcal{M}$ , and we simulate  $\mathcal{M}$  until it gets close to  $\pi$ , at which point the current state of  $\mathcal{M}$  can be considered as a sample from  $\pi$ .  $\mathcal{M}$  is detailed in Algorithm 1. At each state transition,  $\mathcal{M}$  picks a candidate next state  $C$  independently of the current state  $S$  (in Line 6-7). In Line 8, the candidate is either accepted (and  $\mathcal{M}$  moves to  $C$ ) or rejected with certain probability (in which case the candidate state is discarded, and  $\mathcal{M}$  stays at  $S$ ). The main idea is that, at each state, we use the fast but biased sampling technique, which is described in Section 3.1, to propose a candidate  $C$  (in Line 6-7). We correct this bias by adjusting the acceptance/rejection probability (in Line 8) accord-

ingly;  $\mathcal{M}$  is more likely to accept such  $K$ -apps which are less likely to be proposed in Line 6-7. Indeed, as  $\pi(S) = \pi(C)$ , the probability of acceptance is  $\min\left(1, \frac{\Pr[S \text{ is proposed}]}{\Pr[C \text{ is proposed}]}\right) = \min\left(1, \frac{\sum_{u: U_u \supseteq S} 1/\binom{|U_u|}{K}}{\sum_{u: U_u \supseteq C} 1/\binom{|U_u|}{K}}\right) = \min(1, q(S)/q(C))$ . A more formal analysis is described in Appendix A.

---

#### Algorithm 1 MCMC sampling ( $\mathcal{M}$ )

---

```

1: Input: Dataset  $D$ ,  $K$ , # of iterations  $t$ 
2: Output: A sample  $S \in \Omega^K$ 
3: Let  $U := \{D_u : |D_u| \geq K \wedge D_u \in D\}$ 
4: Let  $S$  be an arbitrary  $K$ -apps in  $\Omega^K$ 
5: for  $k = 1$  to  $t$  do
6:   Select an individual  $u \in [1, |U|]$  uniformly at random
7:   Select a subset  $C \subseteq U_u$  uniformly at random such that
      $|C| = K$ 
8:   Let  $S := C$  with probability  $\min(1, q(S)/q(C))$ , where
      $q(x) = \sum_{u: D_u \supseteq x} \prod_{i=1}^K \frac{1}{|U_u| - K + i}$ 
9: return  $S$ 

```

---

**Theorem 1**  $\mathcal{M}$  in Algorithm 1 is an ergodic Markov chain whose unique stationary distribution is the uniform distribution over  $\Omega^K$  for any  $K$ .

The proofs of all theorems in this paper can be found in Appendix A.

**Convergence of  $\mathcal{M}$ .** How to adjust  $t$  in Algorithm 1? We prove that a “good” uniform sample from  $\Omega^K$  can be obtained roughly after  $O(K|D|)$  iterations in most practical cases. The time that  $\mathcal{M}$  takes to converge to its stationary distribution  $\pi$  is known as the *mixing time* of  $\mathcal{M}$ , and is measured in terms of the total variation distance between the distribution at time  $t$  and  $\pi$ .

**Definition 2 (Mixing time)** For  $\xi > 0$ , the mixing time  $\tau_{\mathcal{M}}(\xi)$  of Markov chain  $\mathcal{M}$  is

$$\tau_{\mathcal{M}}(\xi) = \min\{t' : \|P_{\mathcal{M}}^{t'} - \pi\|_{tv} \leq \xi, \forall t \geq t'\}$$

where  $\|P_{\mathcal{M}}^t - \pi\|_{tv} = \max_{x \in \Omega^K} \frac{1}{2} \sum_{y \in \Omega^K} |P_{\mathcal{M}}^t(x, y) - \pi(y)|$  defines the total variation distance.  $P_{\mathcal{M}}^t(x, y)$  denote the  $t$ -step probability of going from state  $x$  to  $y$ , and  $P_{\mathcal{M}}^t$  denote the  $t$ -step probability distribution over all states.

The next theorem shows that  $\mathcal{M}$ ’s mixing time is  $O(|D| \log(1/\xi)/H_1^*)$ , where  $|D|$  is the dataset size and  $H_1^*$  is the unicity of  $K$ -apps from the largest record of  $D$ . As the unicity of  $K$ -apps is usually large in practice, especially if  $K$  is large,  $\mathcal{M}$  is fast-mixing in general. In our dataset  $D$ ,  $0.6 \leq H_1^* \leq 0.999$  for  $2 \leq K \leq 9$ <sup>3</sup>.

**Theorem 2 (Mixing time of  $\mathcal{M}$ )** Let  $H_1^*$  denote the probability that a randomly selected set of  $K$  items from the largest record (i.e., having the most apps) in  $D$  is unique. Then,  $\tau_{\mathcal{M}}(\xi) \leq |D| \ln(1/\xi)/H_1^*$  for any  $K$ .

<sup>3</sup>The unicity of  $K$ -apps from a single record can easily be approximated with Inequality 2 using uniform samples over all  $K$ -apps from the record. Likewise the biased sampling in Section 3.1, this sampling is easy to implement (e.g., by choosing  $K$  items individually from the record without replacement).

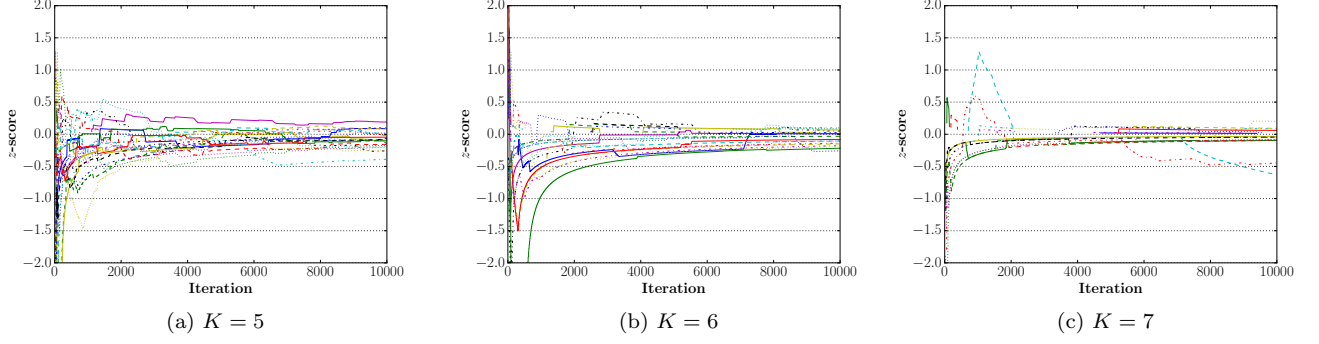


Figure 1: Convergence of our Markov chain  $\mathcal{M}$ . The  $z$ -score, depending on the number of iterations  $t$ , of 20 independent chains are plotted.

We emphasize that the bound in Theorem 2 is a worst-case bound, and the real convergence time can be much smaller depending on the dataset  $D$  as well as the starting state of the chain. As we show next,  $\mathcal{M}$  indeed exhibits much smaller convergence time than its theoretical worst-case bound for our dataset. We detected the convergence of  $\mathcal{M}$  using the Geweke diagnostic [5]; if  $X_t$  denotes a Bernoulli random variable describing whether the current state of  $\mathcal{M}$  at time  $t$  is unique, and  $\mathbf{X}_t = (X_1, X_2, \dots, X_t)$ , then we compute the  $z$ -score  $z = \frac{E[\mathbf{X}_a] - E[\mathbf{X}_b]}{\sqrt{\text{Var}(\mathbf{X}_a) + \text{Var}(\mathbf{X}_b)}}$ , where  $\mathbf{X}_a$  is the prefix of  $\mathbf{X}_t$  (first 10%), and  $\mathbf{X}_b$  is the suffix of  $\mathbf{X}_t$  (last 50%). We declare convergence when the  $z$ -score falls within  $[-1, 1]$ . Indeed, if  $\mathbf{X}_a$  and  $\mathbf{X}_b$  become identically distributed (i.e.,  $\mathbf{X}_a$  and  $\mathbf{X}_b$  appear to be uncorrelated), the  $z$  values become normally distributed with mean 0 and variance 1 according to the law of large numbers. We simulated 20 instances of  $\mathcal{M}$  each starting at different states, and plotted the  $z$ -score of each chain depending on the number of iterations  $t$  in Figure 1. This shows that convergence is detected roughly after 3000 steps in all chains with different values of  $K$ . When this happens, the current state can be taken as a valid sample. Hence, in the sequel, we run  $\mathcal{M}$  with  $t = 3000$  to obtain a uniform sample from  $\Omega^K$ .

We note that  $q$  in Algorithm 1 can be computed rapidly in practice by precomputing another dataset  $T$ , where each record corresponds to an application in  $D$ , and record  $i$  contains the sorted list of all users who have application  $i$  in their record. Hence, the set of users who have a common specific  $K$ -apps can be computed easily by taking the intersection of the corresponding records in  $T$ . The complexity of this operation is  $O(K|i_{\max}|)$ , where  $|i_{\max}|$  is the maximum record size in  $T$ , i.e., the number of users of the most popular application in  $D$ . Fast implementations of the intersection of sorted integers are described in [7].

### 3.3 Computing the sample size

In order to compute the sample size, we use the Chernoff-Hoeffding inequality [6] on the tail distribution of the sum of independent (but not necessarily identically distributed) Bernoulli random variables. In particular, if  $X_i$  denotes a Bernoulli random variable describing the event that the  $i$ th sampled  $K$ -apps is unique in  $D$ , then the deviation of the estimator  $\hat{H}_1 = \sum_{i=1}^n X_i/n$  from  $E[\hat{H}_1] = H_1$  is given by

$$\Pr \left[ \left| \hat{H}_1 - H_1 \right| \geq \varepsilon \right] \leq 2e^{-2n\varepsilon^2}, \text{ or equivalently, } \Pr \left[ \left| \hat{H}_1 - H_1 \right| < \varepsilon \right] \geq 1 - 2e^{-2n\varepsilon^2} \quad (1)$$

where  $\varepsilon$  is the sampling error and  $\sigma = 1 - 2e^{-2n\varepsilon^2}$  is the confidence. Hence, we obtain that

$$n \geq \frac{1}{2\varepsilon^2} \ln \left( \frac{2}{1 - \sigma} \right) \quad (2)$$

For instance, for  $\varepsilon = 0.01$  and  $\sigma = 0.99$ , we need to sample at least 26492  $K$ -apps from  $D$  (with replacement). This guarantees that  $|\hat{H}_1 - H_1| < 0.01$  with probability at least 0.99.

Considering RAD, suppose we aim at approximating the first  $k$  relative frequency values of  $\mathbf{H}$ , i.e.,  $(H_1, H_2, \dots, H_k)$ . Therefore, we wish to simultaneously satisfy Inequality 1 for each  $H_i$  ( $1 \leq i \leq k$ ), where  $\hat{H}_i = \sum_{j=1}^n X'_j/n$ , and  $X'_j = 1$  if the  $j$ th sampled  $K$ -apps occurs in exactly  $i$  records of  $D$ , otherwise  $X'_j = 0$ . Hence,

$$\Pr \left[ \bigwedge_{i=1}^k \left| \hat{H}_i - H_i \right| < \varepsilon \right] \geq 1 - \sum_{i=1}^k \Pr \left[ \left| \hat{H}_i - H_i \right| \geq \varepsilon \right] \geq 1 - 2ke^{-2n\varepsilon^2}$$

where  $\delta = 1 - 2ke^{-2n\varepsilon^2}$  is the confidence. Therefore,

$$n \geq \frac{1}{2\varepsilon^2} \ln \left( \frac{2k}{1 - \sigma} \right) \quad (3)$$

For instance, for  $\varepsilon = 0.01$ ,  $\sigma = 0.99$ , and  $k = 10$ , we need to sample at least 38005  $K$ -apps from  $D$  (with replacement). This will guarantee that  $|\hat{H}_i - H_i| < 0.01$  for all  $1 \leq i \leq k$  with probability at least 0.99.

## 4. EVALUATION

### 4.1 Dataset characteristics

The analyzed dataset comes from the Carat research project [12]<sup>4</sup>. The goal of the project is to inform users about the energy consumption profile of their running apps and give personalized recommendations for improving the battery life. To do so, users need to install the Carat app on their smartphones which intermittently takes measurements about the

<sup>4</sup><http://carat.cs.helsinki.fi>

device, especially when the battery level changes. One type of information collected (among others) by the Carat app is the list of running apps on the device.

The Carat app is available on both Android and iOS. However, our dataset contains only Android users of the Carat app. The dataset includes data from 54,893 Carat users between March 11, 2013 and October 15, 2013 [15]. During this period, the Carat app was collecting the list of running apps (and not the list of all installed apps) on the device when the battery level changes. *As collecting the list of running apps multiple times over more than 7 months is likely to sum up to the set of all installed apps of a user, we consider a record as the set of installed applications in this paper, even if a record might not be the complete set of installed apps all the time.* Also, as system apps are common to all users, we do not consider them for our study, and therefore removed them from all records.

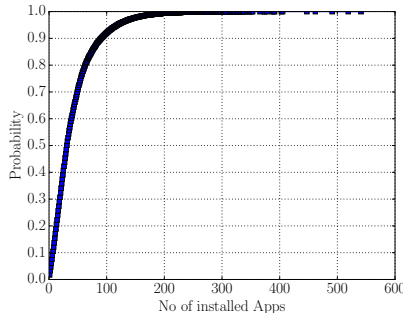


Figure 2: Cumulative distribution of the number of installed apps

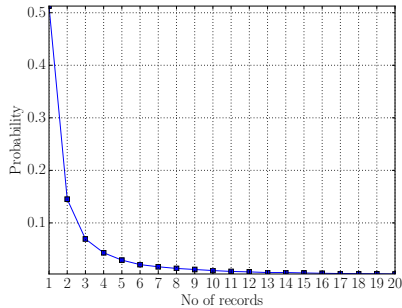


Figure 3: Probability distribution of the number of records containing a particular app

Without system apps, our analyzed dataset contains 92,210 different applications whereas the total number of apps available on the GooglePlay was around 1 million during this time<sup>5</sup>. Furthermore, the average number of apps installed per user in our dataset is 42 with a standard deviation of 39. Table 1 summarizes the main characteristics of our dataset  $D$ .

Figure 2 depicts the cumulative distribution of the number of apps installed by a particular user. We note that more than 90% of users have 100 or fewer applications. Probability

<sup>5</sup>[http://en.wikipedia.org/wiki/Google\\_Play](http://en.wikipedia.org/wiki/Google_Play)

Dataset size $ D $	54,893
# of all apps in $D$	92,210
Maximum record size $\max_u  D_u $	541
Minimum record size $\min_u  D_u $	1
Average record size	42
Std.dev of record size	39

Table 1: Characteristics of our dataset  $D$

distribution of the number of users who installed a particular app is depicted by Figure 3. Notice that more than half of the apps are contained by only a single record in  $D$ .

**Ethical Considerations.** The analyzed dataset comes from the Carat research project [12]. The data were collected with the users' consent, and they were explicitly informed that their data could be used and shared for various research projects. In fact, the Carat privacy policy (available at <http://carat.cs.helsinki.fi>) clearly specifies that "Carat is a research project, so we reserve the right to publish our results online and in academic publications. We also reserve the right to release the data sets into the public domain."

The dataset was shared with us by the Carat team in a pseudo-anonymised form. In particular, identifiers were removed, and each application name was replaced with its SHA1 hash. It contained 54,893 records [15], i.e. one record per user. Each record is composed of the list of applications installed by the user. Furthermore, the data sharing agreement that we signed, stipulated that we cannot use the dataset to deanonymize the users in the dataset.

## 4.2 Results

We find that 98.93% of users have unique set of installed apps in  $D$ , i.e., there does not exist any other user with the same set of installed apps. This means that if we know the list of all the installed apps of a user in the dataset, we can identify that user in the dataset with a probability of 0.99. As the adversary might not always be aware of all the installed apps of a user in practice, we measure the unicity of  $K$ -apps for different values of  $K$  using our dataset  $D$ .

Figure 4 gives the unicity of  $K$ -apps with different values of  $K$  (changing from 1 to 10) for the two different types of sampling techniques described in Section 3: the biased sampling from [2, 3] and our unbiased, uniform sampling described in Section 3.2. In each case, we computed the sample size using Inequality 1 with maximum sampling error  $\varepsilon = 0.01$  and confidence  $\sigma = 0.99$ . Otherwise stated explicitly, we use this sample size in the sequel. This results in 26492 samples for each value of  $K$ . As biased sampling favours more popular  $K$ -apps, the sample unicity  $\hat{H}_1$  is less than with our unbiased approach. In particular, the difference can be as large as 0.5 for smaller values of  $K$ , while it decreases as  $K$  increases. For the unbiased estimation, the sample unicity is 0.75 with  $K = 2$ , and it reaches 0.99 when  $K = 6$ .

Figure 4 shows that the unicity of any  $K$ -apps is large and hence there would be a real privacy threat if such dataset was released. Moreover, Figure 5 depicts the relative abun-

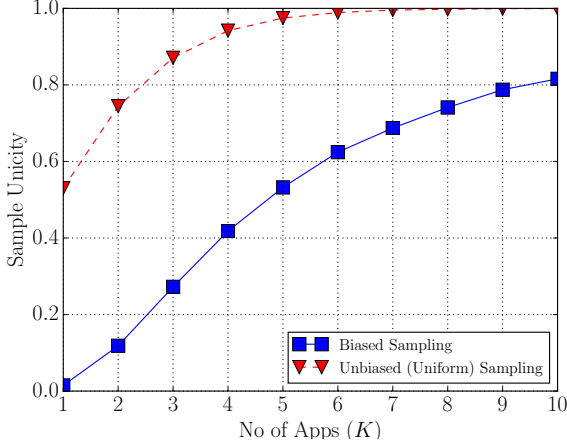


Figure 4: Uniqueness probability as a function of  $K$  for biased and unbiased sampling

dance distribution in  $D$ , when  $1 \leq K \leq 8$ . RAD provides complementary information about users' privacy in  $D$ . In particular, even if the adversary cannot single out the record of the target user in  $D$ , it might still learn new information about him/her. For example, if the known  $K$ -apps of the target user are shared by multiple users in  $D$  and all these users have some identical apps besides the known  $K$ -apps, then the adversary learns that the target user also has these apps installed on his/her phone. This attack is often referred to as the homogeneity attack in the literature [9]. We computed the required sample size using Inequality 3 with  $\varepsilon = 0.01$  and  $\sigma = 0.99$  for  $k = 20$ . This gives 41470 samples overall, which were taken with our uniform sampler  $\mathcal{M}$ .

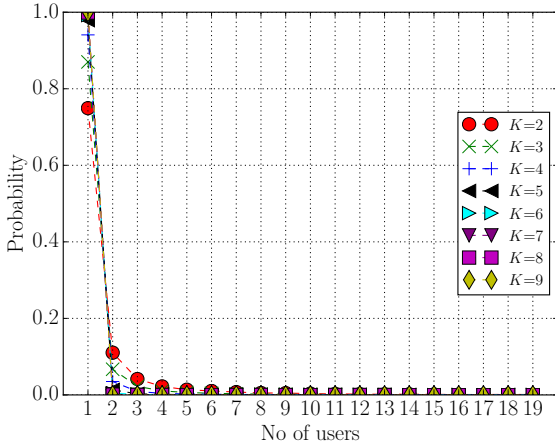


Figure 5: Relative abundance distribution of apps for different sizes of sets of apps

To study the effect of number of users on unicity, we randomly select subsets of users of different sizes from  $D$ , and calculate the sample unicity within these subsets. Figure 6 depicts how unicity changes with the number of users in our dataset. We find that unicity decreases if the user number increases. However, this decrease becomes less significant

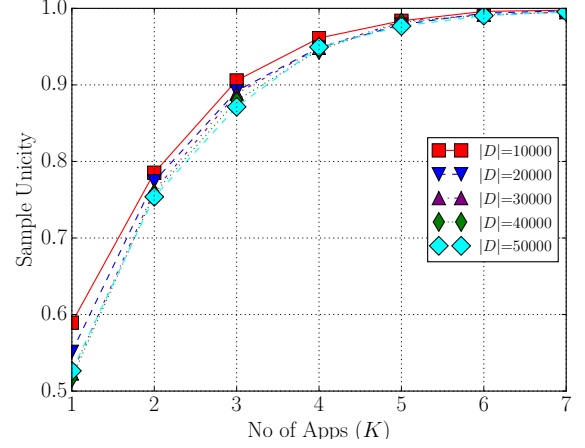


Figure 6: Effect of number of users

for larger number of users. This is probably due to the fact that the number of apps starts to saturate if the user number increases.

As the size of our dataset is much less than the population size of all Android users worldwide (which was roughly 1 billion as of 2014<sup>6</sup> with 1.2 million different applications available on GooglePlay<sup>7</sup>), we aim at predicting the unicity in a larger dataset (possibly in the whole population) in the next section.

## 5. UNICITY GENERALIZATION FOR LARGER DATASETS

Information surprisal can be used to measure uniqueness in the population  $D$  [4]. In our case, the population is all the Android users worldwide, to which we want to generalize our results. As information surprisal of any  $K$ -apps  $\{A_1, A_2, \dots, A_K\}$  over  $D$  is equal to  $-\log(\Pr[A_1, A_2, \dots, A_K])$ , we must need to first measure the co-occurrence probability  $\Pr[A_1, A_2, \dots, A_K]$  of these apps in  $D$ . The co-occurrence probability can be easily computed if we can assume that apps co-occur independently in the dataset as  $\Pr[A_1, A_2, \dots, A_K] \approx \prod_{i=1}^K \Pr[A_i]$ , and  $\Pr[A_i]$  (the popularity of app  $A_i$ ) can be obtained from the download count of  $A_i$  available on Google PlayStore. However, this is not the case in a real-world scenario as there exist correlation between apps installed by a user. As our dataset is very likely to be too limited to capture this correlation (e.g., our dataset contain only 93K distinct apps whereas there are more than 1.2 million available apps on GooglePlay), we cannot take this approach to measure the uniqueness in the population of Android users. We rather employ regression analysis on our dataset which does not rely on this correlation information to predict the unicity in a larger dataset.

For regression analysis, we randomly create datasets of dif-

<sup>6</sup><http://www.engadget.com/2014/06/25/google-io-2014-by-the-numbers/>

<sup>7</sup><http://www.appbrain.com/stats/number-of-android-apps>



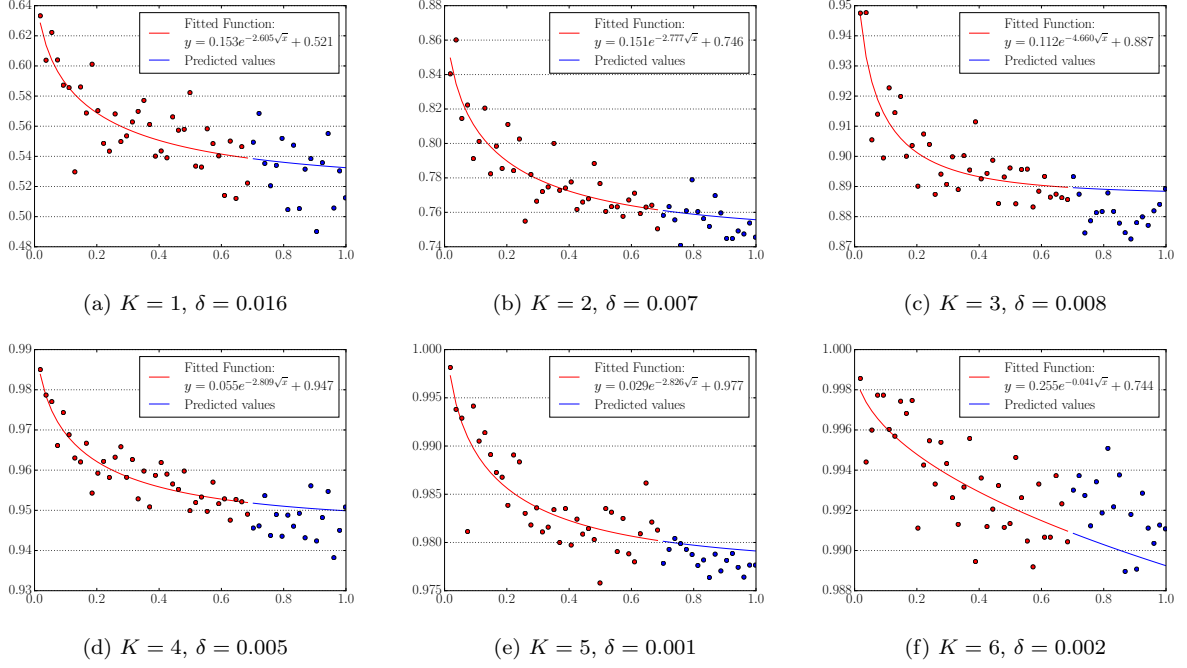


Figure 7: Unicity generalization for different values of  $K$ , trained all with maximum 37000 users. The learnt models (i.e.,  $f(x)$ ) are present in the legend.  $x$ -axis corresponds to normalized dataset sizes with a normalization factor of  $1/54893$ , and  $y$ -axis depicts sample unicity.

ferent sizes from our original dataset and compute the sample unicity for these datasets of different sizes. This gives us the tuples  $(x, y)$  where  $x$  is the number of users in a particular dataset (independent variable) and  $y$  is the calculated unicity value dependent on  $x$ . Here, we assume that the unicity value  $y$  only depends on the number of users  $x$ . We must note that, in reality, unicity depends on many factors such as the characteristics of the users, how many (un)popular applications users tend to have, etc. As it is difficult to take into account all these factors either because they are unknown or hard to measure, we assume that unicity in general is a “proper” function of only the total number of users in the dataset. That is, all other dependent factors are implicitly incorporated into the model, i.e., the general form of the function.

Once we have these  $(x, y)$  tuples, our goal is to select the best model and its parameters that capture the relation between  $x$  and  $y$ . The overall approach is as follows: we divide our  $(x, y)$  tuples in training and test sets. We select the best model (i.e., a function family) based on the general characteristics of application unicity and then learn its exact parameters using our training set. Finally, using the best model thus obtained, we test its accuracy on the test set. This model should be able to predict the unicity value for any dataset of arbitrary size.

**Training and Testing.** We divided our original dataset of 54,893 users in 54 smaller datasets, each of size varying from 1k to 54k. We take the first 70% of all  $(x, y)$  points for training and the last 30% (corresponding to larger datasets) for testing. We deliberately take the last points corresponding

to larger datasets for testing set because we aim to evaluate our model performance on larger datasets, i.e., we want to test how accurately the learned model could be extrapolated.

As we divided our datasets by randomly selecting users out of the original dataset, users in the training and testing set may overlap. However, we found that unicity merely depends on the number of users in the dataset and not specifically on the underlying individuals. For example, we computed the unicity of 50 different sets of 1000 users selected randomly, and found out that the variance of the measured sample unicity is very small.

**Model selection.** To select our model, we first tried linear regression with non-linear basis functions (polynomials of various orders) with and without regularization. However, they provided very inaccurate predictions of unicity. Finally, we selected the following exponential model describing an exponential decay of unicity:

$$f(x) = a \cdot \exp(-b\sqrt{x}) + c \quad (4)$$

The rationale behind choosing this model is as follows. Figure 8 shows that if additional users were added to our dataset, the number of apps would reach the maximum number of apps in the population early as there are fewer apps on GooglePlay than total number of Android users. This suggests that, after a certain point, additional users would not bring many new apps but still, they would bring new combinations of already existing apps. The addition of new combinations of apps should lead to the increase in unicity.



However, the newly added users can lead to the decrease in unicity as well due to the fact that they can also have many already existing combinations of apps. As these two effects of adding new users to the dataset run opposite to each other, we suppose that unicity converges to a value greater than zero which is denoted by  $c$  in Equation 4. Indeed, as Figure 6 shows, although unicity decreases with the increase in the user number, the amount of this decrease tends to decrease as well. A similar observation was made in [2]. We used square root of  $x$  in the exponent in Equation 4 because taking square root is variance-stabilizing<sup>8</sup>. In fact, we also tried other transformations in the exponent but square root yielded the best performance.

The goal of the regression is to compute parameters  $a, b$  and  $c$  in Equation 4 from the training set  $(x, y)$  tuples. In fact, these parameters might be computed employing either standard non-linear regression directly or by first transforming Equation 4 into linear form and then applying linear regression. We use standard non-linear regression because it explicitly computes the lower bound on the unicity value (i.e.,  $c$  in Formula 4). The value of  $x$  is normalized<sup>9</sup> by dividing  $x$  with the maximum size of the dataset for which we want to predict the unicity value.

**Results.** As an error metric, we measure the average absolute error denoted by  $\delta$ , i.e.,

$$\delta = (1/n) \sum_i^n |y_i - f(x_i)|$$

where  $n$  is the number of predicted points,  $y_i$  is the real unicity value, and  $f(x_i)$  is the predicted value.

Figure 7 presents how our exponential model performs on the test set for different values of  $K$ . Although our model can predict the trend of the unicity as the number of users increases in the test set, it slightly overestimates the real unicity in the test set. Nevertheless, the average error  $\delta$  on the test set is only around 0.01. As our app dataset is very small as compared to the whole Android population, we cannot evaluate performance of our model for large number of users, e.g., a few million, or the whole Android population. Therefore, we cannot claim that our model will be able to accurately predict the unicity for datasets having large number of users even if it performs reasonably well on our test data.

**Model validation on a different dataset.** To further demonstrate that our model is a meaningful approach to predict unicity in large populations, we test it on a large mobility dataset provided by a telecom operator in Europe. This dataset contains the Call Data Records (CDR) of 1 million users from a large European city over 6 weeks. Each record in the dataset corresponds to a user and contains the set of his/her visited cell towers, where the total number of different cell towers is 1303. From this dataset, we created smaller datasets of different sizes ( $x$  ranging from 1000 to

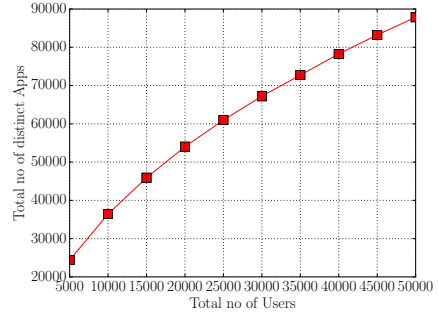


Figure 8: No of distinct apps installed by users

1 million users). Then, we trained our model on the first 6 points (i.e., until the dataset size of 50,000). Figure 9 shows that the model does not predict accurately the unicity for large datasets and the error can be as large as 0.6. Next, we trained the model on the first 7 points (i.e., until the dataset size of 75,000). In this case, we find that the model performs significantly better than in the previous case with an error of 0.13 on average. Finally, we trained the model on the first 8 points (i.e., until the dataset size of 100,000). In this case, we find that the model have accurate predictions for larger datasets, e.g., for a test dataset of size 1 million (10 times more than the maximum size of the dataset used in the training phase), and the error is 0.05 on average.

We find that a mobility dataset of 0.1 million users is sufficient to learn an accurate model and predict the unicity values for a larger mobility dataset. However, as we saw earlier, the model is not able to predict well the unicity of a large population if it is trained on a dataset of only 50,000 users. This may suggest that 50,000 users might be too small in general to learn an accurate model and therefore, our app dataset of 50,000 users might not be sufficient to learn the model. On the other hand, even if this model performs well on a mobility dataset with 1 million users, it does not necessarily imply its good performance on large application datasets due to the different data and user characteristics. Nevertheless, these two experiments together show that our exponential model can be a meaningful approach to predict unicity in large populations.

## 6. CONCLUSION

The paper shows that the list of installed applications by users is quite unique. This result has few implications on user's privacy. First, since this metadata is unique, it could easily be used to profile users based on, for example, the category of the apps. This is what Twitter is doing to provide interest-based targeted ads to users. Second, as a combination of even small number of installed apps is quite unique, this information could be used to re-identify users in a dataset. For example, if Twitter decided to publish the list of apps installed by its users on their smartphones, it would be easy for anyone, who knows 4 or 5 apps of a given user, to re-identify him and discover other apps that are also installed on his smartphone. This makes anonymization of this information challenging, and this is part of our future work.

<sup>8</sup>[https://en.wikipedia.org/wiki/Variance-stabilizing\\_transformation](https://en.wikipedia.org/wiki/Variance-stabilizing_transformation)

<sup>9</sup>[https://en.wikipedia.org/wiki/Feature\\_scaling](https://en.wikipedia.org/wiki/Feature_scaling)

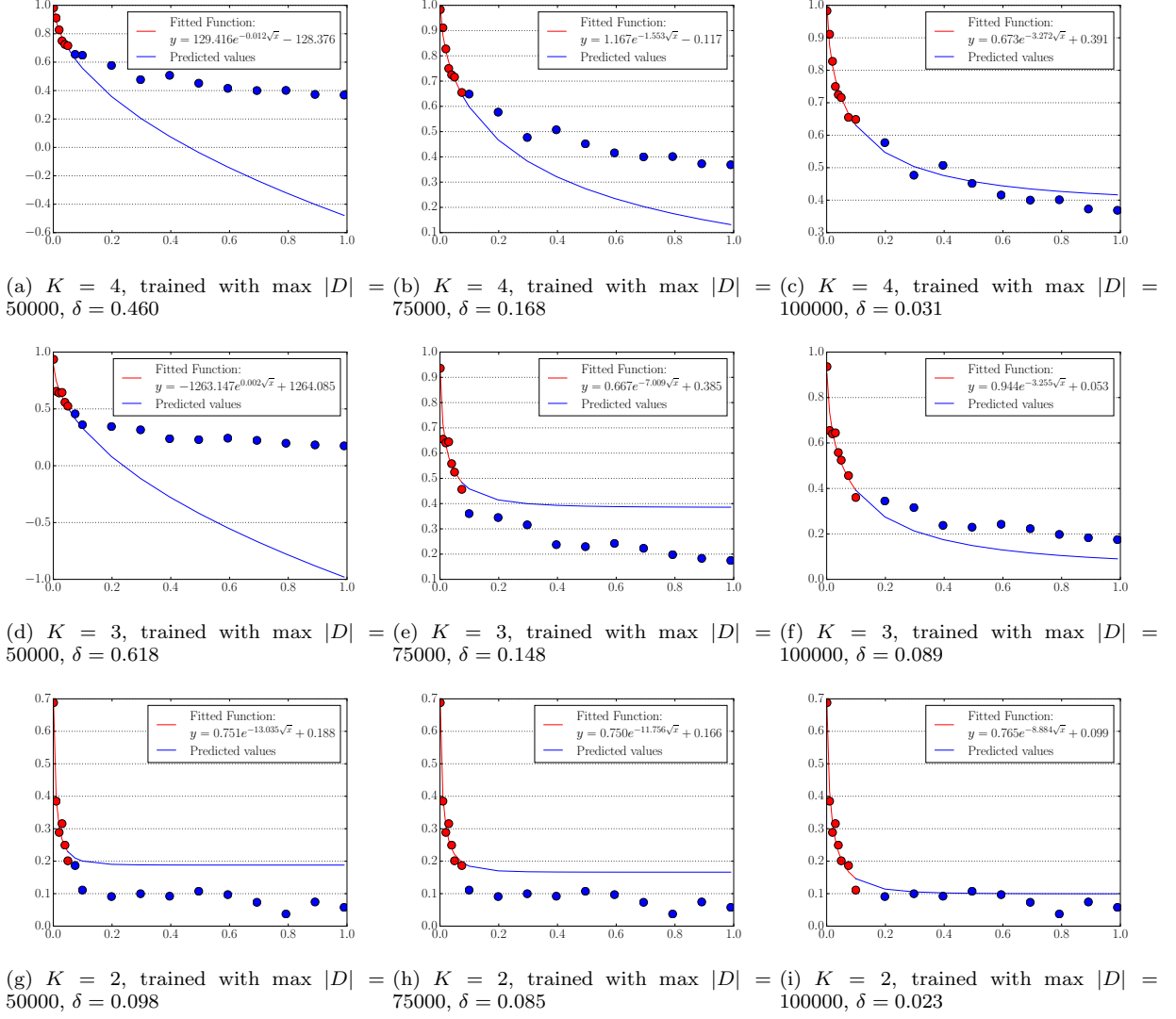


Figure 9: Unicity generalization for different values of  $K$  for location data, trained with datasets of different sizes. The learnt models (i.e.,  $f(x)$ ) are present in the legend.  $x$ -axis corresponds to normalized dataset sizes with a normalization factor of  $1/10^6$ , and  $y$ -axis shows sample unicity.

In general, mobile users reveal many pieces of information that, when combined together, provide a lot of information about users and can be used to build personalized profiles. Since people are unique in many different known and unknown ways, preserving the privacy of mobile users is very challenging. New protection measures need to be devised.

## Acknowledgements

We thank the whole Carat team for sharing their dataset with us. Also, our special thanks go to H. Truong, N. Asokan and S. Tarkoma for discussions related to the dataset. This work is partially funded by Inria project lab CAPPRIS.

## 7. REFERENCES

- [1] S. Chib and E. Greenberg. Understanding the metropolis-hastings algorithm, 1995.
- [2] Y.-A. de Montjoye, C. A. Hidalgo, M. Verleysen, and V. D. Blondel. Unique in the crowd: The privacy bounds of human mobility. *Scientific Reports, Nature*, March 2013.
- [3] Y.-A. de Montjoye, L. Radaelli, V. K. Singh, and A. Pentland. Unique in the shopping mall: On the reidentifiability of credit card metadata. *Science*, 347(6221), January 2015.
- [4] P. Eckersley. How unique is your web browser? In *PETS*, 2010.
- [5] J. Geweke. Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. In *In Bayesian Statistics*, pages 169–193. University Press, 1992.
- [6] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [7] D. Lemire, L. Boytsov, and N. Kurz. SIMD compression and the intersection of sorted integers.

CoRR, abs/1401.6399, 2014.

- [8] D. A. Levin, Y. Peres, and E. L. Wilmer. *Markov Chains and Mixing Times*. American Mathematical Society, 2009.
- [9] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam. *L*-diversity: Privacy beyond *k*-anonymity. *TKDD*, 1(1), 2007.
- [10] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, (6):1087–1092.
- [11] L. Olejnik, C. Castelluccia, and A. Janc. On the uniqueness of web browsing history patterns. *Annals of Telecommunications*, 69(1), February 2014.
- [12] A. J. Oliner, A. P. Iyer, I. Stoica, E. Lagerspetz, and S. Tarkoma. Carat: Collaborative energy diagnosis for mobile devices. In *ACM SenSys*, 2013.
- [13] S. Seneviratne, A. Seneviratne, P. Mohapatra, and A. Mahanti. Predicting user traits from a snapshot of apps installed on a smartphone. *SIGMOBILE Mob. Comput. Commun. Rev.*, 18(2), 2014.
- [14] L. A. Sweeney. *Computational Disclosure Control: A Primer on Data Privacy Protection*. PhD thesis, Cambridge, MA, USA, 2001.
- [15] H. T. T. Truong, E. Lagerspetz, P. Nurmi, A. J. Oliner, S. Tarkoma, N. Asokan, and S. Bhattacharya. The company you keep: Mobile malware infection rates and inexpensive risk indicators. In *WWW*, 2014.

## APPENDIX

### A. PROOF OF THEOREM 1

**Metropolis-Hastings algorithm.** Consider an ergodic Markov chain  $\mathcal{G}$  with transition graph  $G(\Lambda, E)$  and transition matrix  $P_{\mathcal{G}}$ , where  $\Lambda$  is the finite state space. For each  $x \in \Lambda$ , let  $\kappa : \Lambda \times \Lambda \rightarrow [0, 1]$  denote a (not necessarily symmetric) proposal probability distribution function such that for all  $x \in \Lambda$ ,  $\kappa(x, x) + \sum_{y \neq x} \kappa(x, y) = 1$ . The transitions of  $\mathcal{G}$  are defined according to the *Metropolis-Hastings rule* as follows. From any state  $x \in \Lambda$ , first select an  $y \in \Lambda$  such that  $(x, y) \in E$  with probability  $\kappa(x, y)$ . Then, “accept” the transition from  $x$  to  $y$  with probability  $\min\left(1, \frac{\pi(y)\kappa(y, x)}{\pi(x)\kappa(x, y)}\right)$ , otherwise stay at  $x$ . It is not difficult to show [8] that such a Markov chain  $\mathcal{G}$  is reversible, i.e.,  $\pi(x)P_{\mathcal{G}}(x, y) = \pi(y)P_{\mathcal{G}}(y, x)$  and therefore its stationary distribution  $\pi$  is unique [8]. Consequently, after sufficiently many transitions, the distribution of states will be very close to  $\pi$ . Notice that there is no need to compute the normalization constant of  $\pi$ , even if  $|\Omega|$  is very large (i.e., an exponential function of  $K$  in our problem), because it appears both in the numerator and denominator of the transition probabilities.

**PROOF OF THEOREM 1.** In each iteration,  $\mathcal{M}$  can select any individual  $u$  in  $D$ . Hence, at any state,  $\mathcal{M}$  can visit any state in  $\Omega^K$ . Therefore,  $\mathcal{M}$  is connected and aperiodic. Also

notice that

$$\begin{aligned} \frac{\pi(C)\kappa(C, S)}{\pi(S)\kappa(S, C)} &= \frac{\kappa(C, S)}{\kappa(S, C)} \\ &= \frac{\sum_{\forall u: U_u \supseteq S} 1/\binom{|U_u|}{K}}{\sum_{\forall u: U_u \supseteq C} 1/\binom{|U_u|}{K}} \\ &= \frac{\sum_{\forall u: U_u \supseteq S} K!/|U| \prod_{i=1}^K \frac{1}{|U_u| - K + i}}{\sum_{\forall u: U_u \supseteq C} K!/|U| \prod_{i=1}^K \frac{1}{|U_u| - K + i}} \\ &= q(S)/q(C) \end{aligned}$$

where  $\pi$  is the uniform distribution over  $\Omega^K$ . Therefore, a candidate next state is accepted with probability  $\min\left(1, \frac{\pi(C)\kappa(C, S)}{\pi(S)\kappa(S, C)}\right) = \min(1, q(S)/q(C))$ , which means that  $\mathcal{M}$  is reversible and its unique stationary distribution is  $\pi$  according to the Metropolis-Hastings rule.  $\square$

### B. PROOF OF THEOREM 2

In order to prove  $\mathcal{M}$ ’s mixing time, we use a standard coupling argument which is described below.

**Definition 3 (Coupling)** A coupling of a Markov chain  $\mathcal{M}$  on state space  $\Omega$  is a Markov chain on  $\Omega \times \Omega$  defining a stochastic process  $(X_t, Y_t)_{t=0}^{\infty}$  such that

- each of the processes  $(X_t, \cdot)$  and  $(\cdot, Y_t)$ , viewed in isolation, is a faithful copy of the Markov chain  $\mathcal{M}$  (given initial states  $X_0 = x$  and  $Y_0 = y$ ); that is,  $\Pr[X_{t+1} = b | X_t = a] = P_{\mathcal{M}}(a, b) = \Pr[Y_{t+1} = b | Y_t = a]$ ; and
- if  $X_t = Y_t$ , then  $X_{t+1} = Y_{t+1}$ .

Condition 1 ensures that each process, viewed in isolation, is just simulating the original chain  $\mathcal{M}$ , and the coupling is designed such that  $X_t$  and  $Y_t$  tend to coalesce (i.e., move closer to each other according to some notion of distance). Once they meet, Condition 2 guarantees that they will move together forward. The time of this coalescence can be used to upper bound the mixing time which is shown by the next lemma.

**Lemma 1 (Coupling lemma [8])** Let  $(X_t, Y_t)_{t=0}^{\infty}$  be a coupling of a Markov chain  $\mathcal{M}$ . For initial states  $x, y$  let  $T^{x,y} = \min\{t : X_t = Y_t | X_0 = x, Y_0 = y\}$  denote the random variable describing the time until  $X_t$  and  $Y_t$  coalesce. Then

$$\|P_{\mathcal{M}}^t - \pi\|_{tv} \leq \max_{x, y \in \Omega} \Pr[T^{x,y} > t]$$

**PROOF OF THEOREM 2.** Define a coupling  $(X_t, Y_t)$  as follows. Let  $X_t$  and  $Y_t$  choose the same individual  $u$  and subset  $C$  in Line 6 and 7 of Algorithm 1, respectively. This is a valid coupling according to Definition 3, since both  $X_t$  and  $Y_t$  are the exact copies of  $\mathcal{M}$ , and they move together after they coalesce.

Let  $p(x) = \kappa(\cdot, x)$  denote the probability that  $x = C$  is selected in Line 6 of Algorithm 1. Let  $X_0 = x$  and  $Y_0 = y$ , and, w.l.o.g.,  $p(x) \leq p(y)$ . Due to the coupling rule,  $X_t$  and  $Y_t$  can coalesce at any time, since  $P_{\mathcal{M}}(x, y) > 0$  for all  $x, y \in \Omega$ . This happens when both  $X_t$  and  $Y_t$  select a state  $z$  such that  $p(z) \leq p(x) \leq p(y)$ , since  $q(z) \leq q(x) \leq q(y)$  will

also hold. Let  $U_{\max} = \max_u U_u$ . For any  $x, z \in \Omega$ , where  $z$  occurs only in  $U_{\max}$ ,  $p(z) \leq p(x)$ . Indeed,

$$\begin{aligned} p(x) &= \frac{1}{|U|} \sum_{\forall u: U_u \supseteq x} \frac{1}{\binom{|U_u|}{K}} \\ &\geq \frac{1}{|U|} \frac{1}{\binom{|U_{\max}|}{K}} \\ &\geq p(z) \end{aligned}$$

Hence,  $X_t$  and  $Y_t$  coalesce as soon as they select any  $z \in \Omega$  which occur only in the largest record in  $D$ . Therefore,

$$\begin{aligned} \|P_{\mathcal{M}}^t - \pi\|_{tv} &\leq \max_{x, y \in \Omega^K} Pr[T^{x, y} > t] \quad (\text{by Lemma 1}) \\ &\leq \sum_{i=t}^{\infty} (1 - H_1^*/|U|)^i H_1^*/|U| \\ &\leq (1 - H_1^*/|U|)^t \\ &\leq \exp(-tH_1^*/|U|) \end{aligned}$$

which proves the theorem.  $\square$