



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Data Sharing in Participatory Social Sensing

Master Thesis

Ramapriya Sridharan

September 3, 2016

Advisors: Prof. Dr. Dirk Helbing, Dr. Pournaras Evangelos
Department of Computational Social Sciences , ETH Zürich

Contents

Contents	i
1 Introduction	3
2 Related Work	5
3 Computational Model	9
3.1 Introduction	9
3.2 Model Intricacies	9
3.2.1 Collecting User Information	9
3.2.2 Categorization of the Features	10
3.2.3 Categorization of the Sub-Features	12
3.2.4 Weight Matrix Calculation	13
3.2.5 Cost Matrix Calculation	14
3.2.6 Cost and Privacy Metrics	14
3.2.7 Improving the Metrics	16
3.2.8 Summarization of Collected Data	16
3.3 Analysis of the Model	17
3.3.1 Setup	17
3.3.2 Results	18
4 Experiment Methodology	25
4.1 Preparatory Phase	25
4.1.1 Pre-Survey	25
4.1.2 Sub-Features	26
4.1.3 Privacy Options	27
4.1.4 Question Structure	28
4.1.5 Budget and Experiment Duration	29
4.2 Entry Phase	30
4.2.1 Collecting General User Information	30

CONTENTS

4.2.2	Categorization of Features	30
4.2.3	Categorization of Sub-Features	31
4.2.4	Answering Questions with No Incentives	32
4.3	Core Phase	33
4.3.1	Improve Privacy or Credit	38
4.3.2	Answering Questions with Incentives	39
4.4	Exit Phase	40
4.5	FairDataShare Web Portal	41
4.5.1	Data Generator's Portal	41
4.5.2	Stakeholder's Portal	41
5	Explanation of the Mobile Application	45
5.1	The Building Blocks	45
5.2	The Mobile Application	45
5.2.1	Local Storage	46
5.2.2	Alarms and Notifications	50
5.2.3	Fetching Data Requests	52
5.2.4	Recording User Choices	52
5.2.5	Sensor Data Collection and Summarization	53
5.2.6	Server Synchronization	53
5.3	The Server	55
5.3.1	Kinvey Data Storage	55
5.3.2	FairDataShare Web Portal	62
6	Pre-Survey and Experiment Findings	65
6.1	Overview of the Pre-Survey Data	65
6.2	Pre-Survey Methodology and Findings	67
6.3	Overview of the Experiment Data	79
6.4	Findings from the Experiment Data	79
7	Conclusion	81
A	Appendix	83
	Bibliography	85

Abstract

Data from citizens needs to be collected and analyzed to create or improve current services in society. Data collected from them, in general, reveals information about their behavior and choices. In addition, it can also reveal sensitive information, that they might not be comfortable with. To preserve the privacy of citizens is where data privacy comes into play. There are various methods to maintain data privacy and different levels of privacy to maintain. The higher the privacy level, the more concealed the data is. Given the choice, citizens would generally choose the highest privacy level. At times, less concealed data is needed while solving problems that need data with less errors. To help citizens reduce the level of privacy of the data when needed, different kinds of incentives can be used, such as monetary incentives. From a fixed budget on the demand side, rewards(incentives) are handed out to citizens to incite them to give less privatized data, yet maintaining a minimum level of privacy. The goal of the Thesis is to understand the social dynamics of privacy and information sharing. Existing data can be used or data can be collected for the purpose of the analysis.

Chapter 1

Introduction

Chapter 2

Related Work

Participatory social sensing is the active participation of users with their mobile phones to form a network that enables the collection and analyzation of large amounts of data. The concept was first introduced by Burke et al in 2006, where they talk about the potential benefits and propose an initial architecture. Collecting data from various sensors is important for Big data analysis and to find answers to complex social questions. Lei Song et al [19] performed an extensive survey on the sensor devices and their applications. It was found that different sensor combinations is the foundation to grasp important information. Giannotti et al [11] propose the *Planetary Nervous System* to collect data from connected sensors and use that data to do big data analysis in a privacy aware fashion.

Some applications of participatory sensing are LiveCompare, TraficSense and CenseMe mobile applications. LiveCompare is introduced by Linda Deng et al [10] where they make use of the widespread availability of mobile phones to find cheap groceries making use of the camera for barcode decoding and location to find the stores. TraficSense by Prashanth Mohan et al [15] is a concept aimed to keep track of traffic on the road with a mixture of traffic and vehicle types. IT collects a variety of sensor data such as the accelerometer and location but not limited to. CenseMe created by Emiliano Miluzzol et al [14] where friends in social networks can share their status in terms of their mood, activity, surrounding and habit. This includes physical and virtual sensors that can capture the online life of a person.

Participatory sensing is needed for a fairer system to trade data due to the fact that many mobile applications take data away from users without their knowledge. Jinyan Zang et al [20] did a study using 110 Android and iOS apps to find the ones that share personal information, behavioural information and location data with third parties. Doing this revealed that it does not require a notification from the application. They also found that paid applications still shared sensitive information to third parties.

2. RELATED WORK

Ashwini Rao et al [18] examined the behavioural profiles formed by Google and Yahoo. Participants were surprised and very concerned that data has been collected from them. Additionally, the profiles formed were found to be in some aspects inaccurate and had excess of information, so much that the profiles didn't seem to be anonymous anymore. Further, a survey was created asking participants questions about the behavioural profiles and was launched on Amazon Turk. Participants didn't find the profiles formed to be accessible enough and they also complained that they wanted to know more about who and where the data will be used. Overall, the impression of participants is that the whole process of collecting data lacked transparency. This shows that there is a need for more privacy and control of data from the user side.

Studies have been done to investigate the relationship of users and their data. Alessandro Acquisti et al [1] created a survey to observe the privacy concerns in e-commerce preferences and masking of location data. They found that users do not make reckless decisions, rather they make decisions based on what information they have, how much they care and what they believe the effect of their actions will be. This leans on the fact that with sufficient information users can make rational choices about the privacy of their data.

Rebecca Balebacco et al [2] study through surveys and an experiment that users do not remember the sensors accessed by each application, shown during installation of the mobile application and proposes to inform users during the use of the application itself before collecting the sensor data. Similarly, Lin Jialiu et al [13] examines through crowdsourcing the perception of users to the data collection from mobile applications. The main takeaway from here is that users felt more comfortable if the purpose of a resource access was stated.

Additionally, studies have been done on assessing sensor data sharing in mobile phones. George Danezis et al [9] did a study to assess how much people value their location data using auction technique. They found that the median bid was 43\$ a period of one month, but this varied a lot on whether the person was a student, the relationship status and their travelling habits. Dan Cverck et al [8] also examines the value of location privacy with over 1200 people and varied demographics. In this case the users were told fake goals in order not to be biased about their data privacy. Contradictions were found to the study [9] about the change in value due to travelling habits, but the median bid was found to be the same. There was also differences in results among the users with different demographics. This goes to show that one incentive does not fit all the users.

Delphine Christine et al [7] do an extension of the study in the paper by George Danezis et al [9]. They try to analyse how various factors can affect

data sharing such as demographics, incentives and spatio-temporal elements vary the importance users have on their data for various sensors. They also dive into other aspects such as the purpose for which data is shared and to whom the data is shared. They found that younger people and people with affiliations with buyers of their data tend to share more information. They also found that users claimed more rewards to corporations. The work by Camp Jean [4] has mentioned that the participants of the surveys may not tell the truth inspite of financial rewards. The later only ensures that the users successfully complete the survey.

Other than the surveys, there are studies done on the mobile phones themselves. Brush et al [3] collect the location data of 32 users for a period of 2 months. Users have five privacy options they can choose from:

- Deleting near home
- Mixing to provide k anonymity
- Randomizing
- Discretizing
- Subsampling

At the end of the two months, users were shown visualizations of their data. The authors mentioned that the user interface was not intuitive and that users might have been biased to the location data due to the experimental setup. Additionally, it was found that users were not consistent with privacy decisions and with whom they shared data. It was concluded that users need to be properly informed about every detail to enable them to make rational choices.

Haksoo Choi et al [5] is a framework that provides sharing of sensor data based on rules along with the possibility of applying obfuscation algorithms. They found that users share data with a purpose and hence the purpose of data sharing should be included in the rules. Additionally, Eiji Hayashi et al [12] with 20 participants examine the sensor data sharing with all or nothing options. It is found that all or nothing options are a poor fit for user preferences and sharing of partial sensor information should also be provided.

Various algorithms can be used to protect the privacy of users while providing various amounts of data sharing possibilities. Pournaras et al [17] propose a scheme where users can share various amounts of data. Users supply data to the data aggregators who buy data. The incentives that are received depend on the accuracy of the data that is shared and the accuracy required by the data aggregators. If the data shared is of lower accuracy then the errors in data processing increase. Similarly, with higher accuracy

2. RELATED WORK

the data processing tasks will give lesser errors. The process of manipulating the accuracy of sensor data is called summarization. The errors in the data can be mitigated if there is a large population of users participating in the data sharing tasks. Summarization can be any algorithm from simple arithmetic functions to clustering algorithms.

In this dissertation, it is attempted to address some the drawbacks of the aforementioned studies. A social experiment is created and an Android application developed to study the relationship between incentives and data sharing for a wide range of sensors. Additional factors that can affect data sharing are also included in the equation such as potential data buyers and the purpose for which the data collected is used. Users are allowed to share data in five levels each corresponding to a summarization level. The user interface has been designed to be intuitive and easy to use. Users are approached to share their data and are incentivized with some credits. The amount of credit received is proportional to the amount of data shared. The amount of credit per data request is dynamically allocated according to the user profile created. High amounts are allocated to data requests that are considered privacy intrusive by the user. To keep the users motivated a participation fee is included, which means users will receive credits for answering questions irrespective of the amount of data shared. Additionally, a survey is deployed before the experiment to gain insight into the perception of users on various features and to fine tune the mobile application details.

Delphine Christine et al [6] conclude their paper on the challenges for the future, the following points have been addressed [16]:

- Including the participants in the privacy equation
- Providing composable privacy solutions
- Trade-offs between privacy, performance and data fidelity
- Making privacy measurable
- Defining standards for privacy research
- Holistic architecture blueprints

Chapter 3

Computational Model

3.1 Introduction

The aim is to create a computational model that is able to collect useful information about the influence of monetary incentives on mobile data sharing. (Quote some studies that have done similar studies with no data incentives). The users are first asked some preliminary questions to form a profile about them. The model proceeds to use the user profiles formed to assign each sensor data request with a maximum achievable credit. The model attempts identifies the data requests where users might not be inclined to share mobile sensor data willingly. These data requests are assigned higher maximum obtainable costs. Similarly, the data requests where the users would want to share more mobile sensor data is assigned a lower maximum obtainable cost. This permits us to see whether incentives do indeed make a difference in mobile sensor data sharing. The model aims to identify the amount of data each user would share for a data request and assign maximum obtainable costs accordingly.

3.2 Model Intricacies

The sections below explain the various building blocks of the computational model. The Figure 3.1 provides an overview of the flow of the model.

3.2.1 Collecting User Information

To begin with the model, each user is asked to enter various non-intrusive personal information. The information collected can consist of but is not limited to :

- Gender
- Year of birth

3. COMPUTATIONAL MODEL

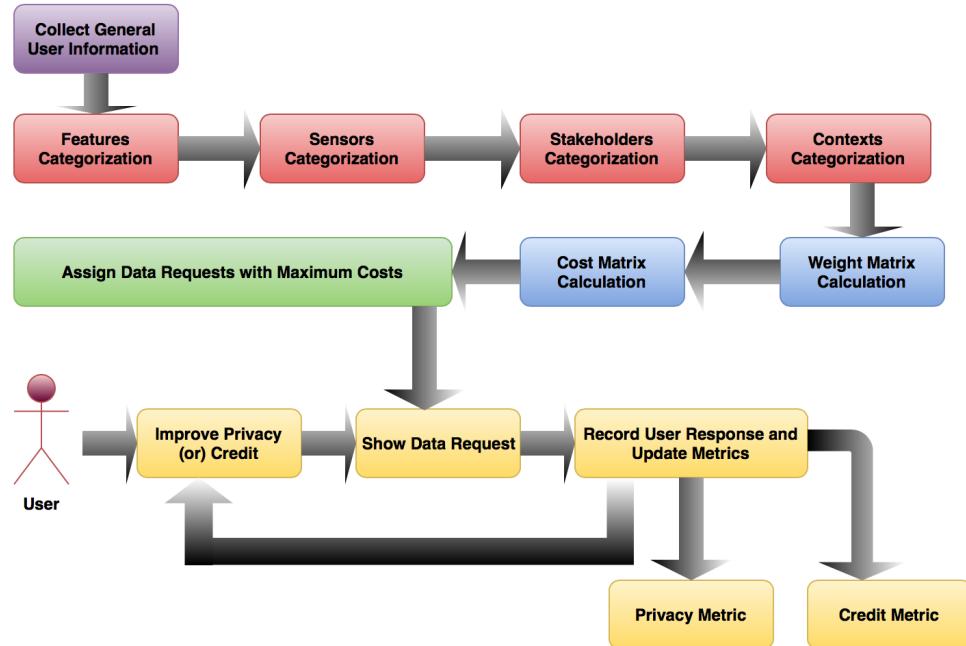


Figure 3.1: Computational Model Flow Chart

- Country
- Education Level
- Occupation
- Frequency of mobile phone use per day
- List of different mobile applications present on the users phones

Collecting this information is crucial to the data analysis that is conducted in the later stages.

3.2.2 Categorization of the Features

After the users personal information is collected, users are asked to place the features in categories according to how privacy intrusive they are to the user. A feature can be one of the following:

- **Sensors** : Sensors consist of the sensors in the mobile phone which users can trade in a data request
- **Stakeholders** : Stakeholders consist of any entity that can request the user for mobile sensor data
- **Contexts** : Contexts consist of the purpose for which a Stakeholder would like to obtain the user's mobile sensor data

3.2. Model Intricacies

Features are the three dimensions that form a unique data request. A data request is defined as a Stakeholder asking users to share their mobile Sensor data for a particular Context.

As mentioned before, users are asked to categorize the features into one of the five categories:

1. Very low privacy intrusion
2. Low privacy intrusion
3. Medium privacy intrusion
4. High privacy intrusion
5. Very high privacy intrusion

Categories are linearly scaled and equally spaced. As indicated by the numbers on the left of the categories, these range from one to five and users can place each of the Features in a category according to their perceived intrusion level. Category one represents that the Feature does not contribute a lot to the data sharing decision. Similarly, category five represents that the Feature contributes a lot to the user's data sharing decision. Similarly, category five represents that users are reluctant to give away their sensor data for this feature. More than one feature can be placed in the same category, which makes it a more powerful tool than the ranking mechanism.

Let the variable cat represent the number of categories, which here is five. Additionally, let the category assigned to the Sensors be represented by the variable se , the category assigned to the Stakeholders be represented by the variable st and the category assigned to the Contexts be represented by the variable co .

Once users have categorized the Sensors, Stakeholders and the Contexts into the respective categories reflecting the importance of each of the features in the data sharing decision, each feature is assigned a weight. Let the respective weights of Sensors, Stakeholders and Contexts be represented by the variables, w_{se} , w_{dc} and w_{co} are calculated as follows :

$$w_{se} = \frac{se}{se + st + co} \quad (3.1)$$

$$w_{st} = \frac{st}{se + st + co} \quad (3.2)$$

$$w_{co} = \frac{co}{se + st + co} \quad (3.3)$$

3. COMPUTATIONAL MODEL

3.2.3 Categorization of the Sub-Features

Once the features have been categorized and their weights calculated as above, sub-features are to be categorized. A sub-feature is defined as one type of a feature. In other words, sub-features are the different types of features that appear during data request to the user. The following are examples of sub-features for each feature :

- Sensors :
 - Accelerometer
 - Battery
 - Gyroscope
- Stakeholders :
 - Corporation
 - Government
 - Educational Institution
- Contexts :
 - Education
 - Navigation
 - Gaming

Each of the above are different kinds or sub-features of the respective Features. For each of the available features, the respective sub-features need to be in turn categorized in a similar fashion to section 3.2.2. The categories are the same as mentioned in the previous section. Let num_{sf} be the number of sub-features each feature has.

The first category indicates that users find the sub-feature would not hinder the data sharing decision. This means that the user would not be worried trading data for a data request involving this sub-feature. The last category indicates that users find this sub-feature would hinder the data sharing decision . This means that users would be reluctant of giving data for a data request involving this sub-feature. As seen in the conceptual diagram is shown in figure 3.2, users place each of the sub-features available for every feature in the given categories.

Let every sub-feature be represented by a unique identifier within its feature. For example, in the list of sub features provided above, accelerometer

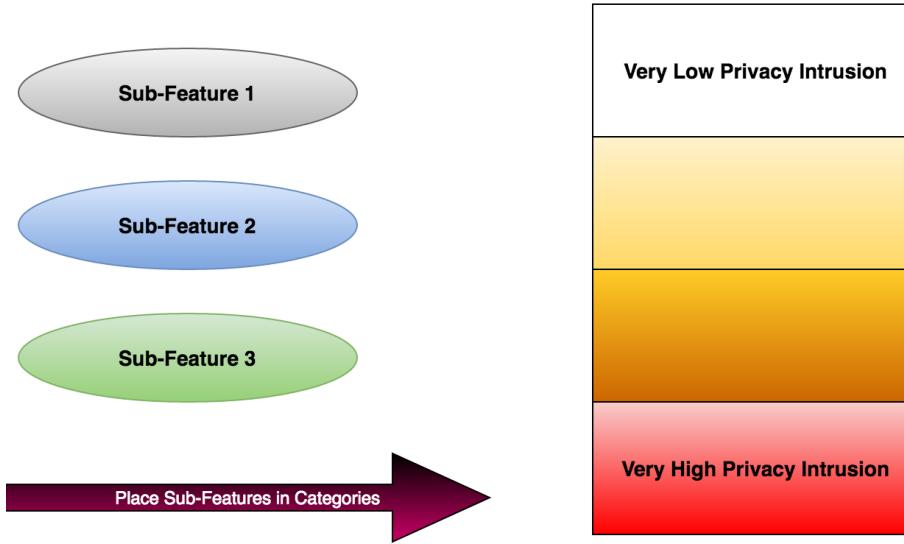


Figure 3.2: Categorizing Sub-Features according to the perceived Intrusion Level

is the first sub-feature of sensors, corporation is the first sub-feature of stakeholders and education is the first sub-feature of contexts. For each of the sub-features of Sensors, categories they are placed in by users is represented by se_i and i is the identifier of the sub-feature. Similarly, categories assigned to sub features of features Stakeholders and Contexts respectively are represented by st_j and co_k , where j and k are the identifiers of the sub-features categorized.

3.2.4 Weight Matrix Calculation

Each data request to users consists of the three above mentioned features in them. Each of the features each have num_{sf} sub-features that can appear in turns in a data request. So the total number of data requests are :

$$num_{dr} = num_{sf} * num_{sf} * num_{sf} \quad (3.4)$$

Let WM be a matrix with three dimensions $num_{sf} \times num_{sf} \times num_{sf}$. We call this the weight matrix. Each cell of WM , that is $WM_{i,j,k}$ represents a data request which involves the Sensors sub-feature with identifier i , Stakeholders sub-feature with identifier j , and the Contexts sub-feature with identifier k . That is, each cell of WM represents the weight of a data request to the users. The aim of the weight matrix is to use the information collected from the user categorizations, to assign weights to each data requests. Intuitively, the process examines the data requests where the user is least likely to trade

3. COMPUTATIONAL MODEL

data and assigns higher weights to those data requests. This process can be seen in section 3.3 with examples. As mentioned before, each cell of the matrix WM represents the weight of a data request with a unique Sensors sub-feature i , Stakeholders sub-feature j and Contexts sub-feature k . To calculate the weight of a data request :

$$WM_{i,j,k} = (se * se_i) + (st * st_j) + (co * co_k) \quad (3.5)$$

Applying this formula for every possible values of i , j and k gives the weight matrix WM .

3.2.5 Cost Matrix Calculation

Once WM has been calculated, it can give an idea of the weight each data request receives. The aim is now to assign a maximum obtainable cost to each data request. This cost is the maximum credit users can receive for a particular data request. Let CM be the cost matrix with the three dimensions $num_{sf} \times num_{sf} \times num_{sf}$. Let it be assumed to have a budget of b for a day, where b can be in an actual currency or any sorts of virtual credits. In this literature the budget will be referred to with the unit credits. Each cell of the cost matrix will represent the amount of credits allocated for a particular data request for one day. To begin with, we calculate the sum of all the cells of the weight matrix WM :

$$sum_{WM} = \sum_{i=1}^{num_{sf}} \sum_{j=1}^{num_{sf}} \sum_{k=1}^{num_{sf}} w_{i,j,k} \quad (3.6)$$

where the function sum_{WM} gives the sum of a matrix, in this case the weight matrix. Let $CM_{i,j,k}$ represent the credit allocated for the data request which involves the Sensor's sub-feature with identifier i , Stakeholder's sub-feature with identifier j , and the Context's sub-feature with identifier k . To calculate one cell of the cost matrix :

$$CM_{i,j,k} = \frac{WM_{i,j,k} * b}{sum_{WM}} \quad (3.7)$$

Repeating the above for every cell of CM , the entire cost matrix can be calculated. Now, all the maximum obtainable costs have been allocated per day for every data request.

3.2.6 Cost and Privacy Metrics

Every data request has been assigned a cost. This is the maximum cost that a user can obtain for that data request. The Cost metric is the total amount

of credits the user has obtained by trading data for data requests for one day. Similarly, the Privacy metric is the amount of privacy percentage obtained while trading data for requests. It intuitively quantifies the amount of data the user has refused to share hence implying privacy. The Cost and Privacy metrics are inversely proportional to each other, in the sense that when the Cost increases the Privacy decreases and vice versa.

Each data request the user chooses how much data is to be shared, from the maximum amount of data to no data at all. The possible responses to a data request are called options. Each option corresponds to a summarization level explained in detail in section 3.2.8. The cost assignment to each option is linearly scaled according to the cost assigned to each data request. Let us assume there are options for a data request ranging from 1 to m (numeric options), where 1 corresponds to the option where the users give all their data for a request and m to where the users choose not give any data for a request. Therefore there are a total of m options for every data request. For each option in a data request is associated with:

- The amount of credit change from this particular option for a data request
- The amount of privacy change from this particular option for a data request

While assigning costs to data requests there are two scenarios to consider:

- Assigning option costs without a participation cost. Users are not rewarded for their participation
- Assigning option costs inclusive of a participation cost. Users are rewarded for responding to data requests irrespective of how much data they choose to share

Let us examine the first scenario. Let the option costs be calculated for the data request with Sensors sub-feature i , stakeholders sub-feature j and contexts sub-feature k . The assigned cost for any option numbered h of this data request is calculated as follows:

$$cost_h = \frac{CM_{i,j,k} * (m - h)}{m - 1} \quad (3.8)$$

Applying this formula by replacing h by the option numbers from 1 to m gives the cost the user can receive for each option.

Similarly, if a participation cost would like to be assigned to each option, it would mean that even tough the user does not share data, they still receive some credit for answering the data request. This concept can be implemented to ensure user participation in the experiment. (Quote some paper

3. COMPUTATIONAL MODEL

with participation of users in PSS). Let x be a fraction of the total budget b that is dedicated for user participation. Using a geometric progression with $a = 1$ and $r = \sqrt[m-1]{x}$, we can calculate the fraction of the maximum cost obtainable from a data request $frac_h$, an option numbered h gets:

$$frac_h = a * r^{h-1} \quad (3.9)$$

The fraction of the cost an option h can be assigned has been calculated, to get the cost $cost_h$ of option h for the data request with Sensors sub-feature i , stakeholders sub-feature j and contexts sub-feature k :

$$cost_h = frac_h * CM_{i,j,k} \quad (3.10)$$

This assigns costs to each option, taking into consideration a participation cost that the user gets even if data is not shared for that data request.

Privacy percentage pri_h is linearly scaled between the first to the m th option between 0% and 100% as follows:

$$pri_h = \frac{(h - 1) * 100}{m - 1} \quad (3.11)$$

The total cost and privacy is the sum and arithmetic average of all the costs and privacy respectively, obtained from every answered data request. If a data request is left unanswered, a maximum privacy of 100% and minimum cost of 0 credit is assumed.

3.2.7 Improving the Metrics

Before users answers a question, it is useful to know what the user interest lies in. Would the user like to improve the privacy metric, or would the user would like to increase the credit revenue. In addition, if we know what the user is looking to improve, we can retrieve the question that can improve the that particular metric the most. For example if the user wishes to improve his privacy further, we look at the questions where the user has given the most amount of data. We then put forth this question to answer, which indicating all the options that can improve the privacy. Similarly, if the user chooses to obtain more credit, the question where the user has given least amount of data is retrieved. Options that can improve the user credit are also indicated.

3.2.8 Summarization of Collected Data

Each data request can have options m number of options the user can choose from for every data request. These options range from 1, which indicates

that the user would like to give all his data, to option number m , which indicates when the user does not want to give any data to this data request. Even though all data is encrypted these days, it is still not enough as encryptions might be cracked. Summarization is a privacy algorithms that aggregates data to provide less information than in its original form. The higher the summarization level gives less data than in its original form. The lower the summarization level gives data closer to its original form. In this model, sensor data is collected for a period of 24 hours every y seconds for every data request. If the data is summarized, according to the option chosen, the data is collected either every y seconds or lesser.

Data is collected for the 24 hour period, and at the end of this period according to the option chosen by the user, it is summarized. Summarization can be linearly assigned to each option. The highest privacy option m corresponds to the highest summarization level. The first option corresponds to the lowest summarization level. An example of assigning the summarization level $summ_h$ for an option h for a data request can be the following :

$$summ_h = y * h \text{ where } h \neq m \quad (3.12)$$

This gives the frequency of sensor data collection for every option of a data request.

3.3 Analysis of the Model

In this section, three different examples are explained in order to test some of the properties that the model assigns to the weight and cost matrix.

3.3.1 Setup

In the following examples, the following features and sub-feature are considered:

1. Sensors
 - a) Accelerometer
 - b) Noise
 - c) Location
2. Stakeholders
 - a) Corporation
 - b) Government

3. COMPUTATIONAL MODEL

- c) Educational Institution
- 3. Contexts
 - a) Navigation
 - b) Environment
 - c) Social Media

The numbers indicated to the left of the sub-features is the sub-feature unique identifier. This uniquely identifies a sub-feature of a feature. There are in total $num_{sf} = 3$ sub-features for each feature. Each user will receive a number of

$$num_{sf} * num_{sf} * num_{sf} = 27$$

data requests in total. The number of categories available to categorize is $cat = 5$ as explained in 3.2.2. Additionally, it is assumed that a budget $b = 100$ per day is available. The input to the model are the user choices during the categorization of the features and sub-features.

3.3.2 Results

To begin with, the way the user has categorized the features and sub-features is introduced. This will be followed by an explanation of the generated matrices. To make reference easier to the graphs, instead of sub-feature names, numeric identifiers are used. From now on each feature and sub-feature will be referred to by its identifier such as feature 1 for Sensors and sub-feature 2 of feature 1 for the noise sensor. The tuple (a,b,c) represents a data request with:

1. a - Sensor's sub-feature a
2. b - Stakeholder's sub-feature b
3. c - Context's sub-feature c

where a, b and c are all numbers from one to three.

Scenario One

If features and sub-features have all been given the same categories respectively by users, then all data requests should be assigned equal weights and costs. In scenario 1, the users choose categories for the Features and sub-features as shown in the table 3.1. As it can be seen in the table, each feature receives the category 1, and all their sub-features are categorized as 3. In short, all the features have the same categorization and their respective sub-features all have the same categorization as well. From this input,

3.3. Analysis of the Model

the formulation of the weight matrix can be seen in figure 3.3, and the cost matrix can be seen in figure 3.4. For each data request indicated as a tuple of (sensors, stakeholders, contexts) in the x-axis of figures ??, all have identical weights and costs. This is due to the fact that the users find all the features and sub-features equally intrusive so all the data requests are weighted equally.

Table 3.1: Categorization for Scenario 1

Feature	Sub-Feature ID = 1	Sub-Feature ID = 2	Sub-Feature ID = 3
Sensors 1	Accelerometer 3	Noise 3	Location 3
Stakeholders 1	Corporation 3	Government 3	Educational Institution 3
Contexts 1	Navigation 3	Environment 3	Social Media 3

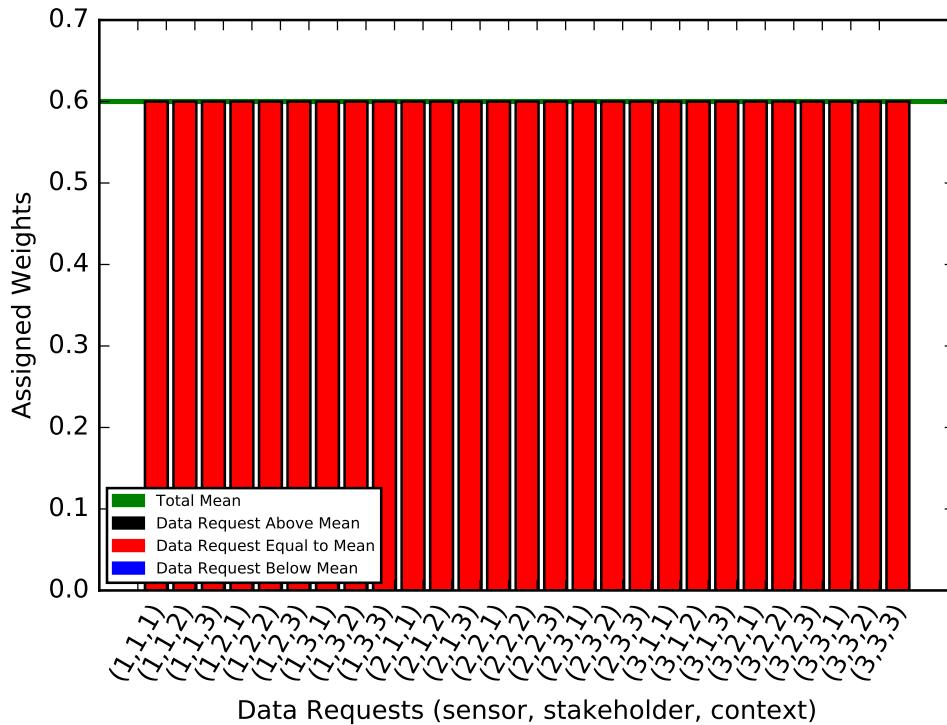


Figure 3.3: Values of the Weight Matrix

We can conclude that if the users perceive the feature and respective sub-

3. COMPUTATIONAL MODEL

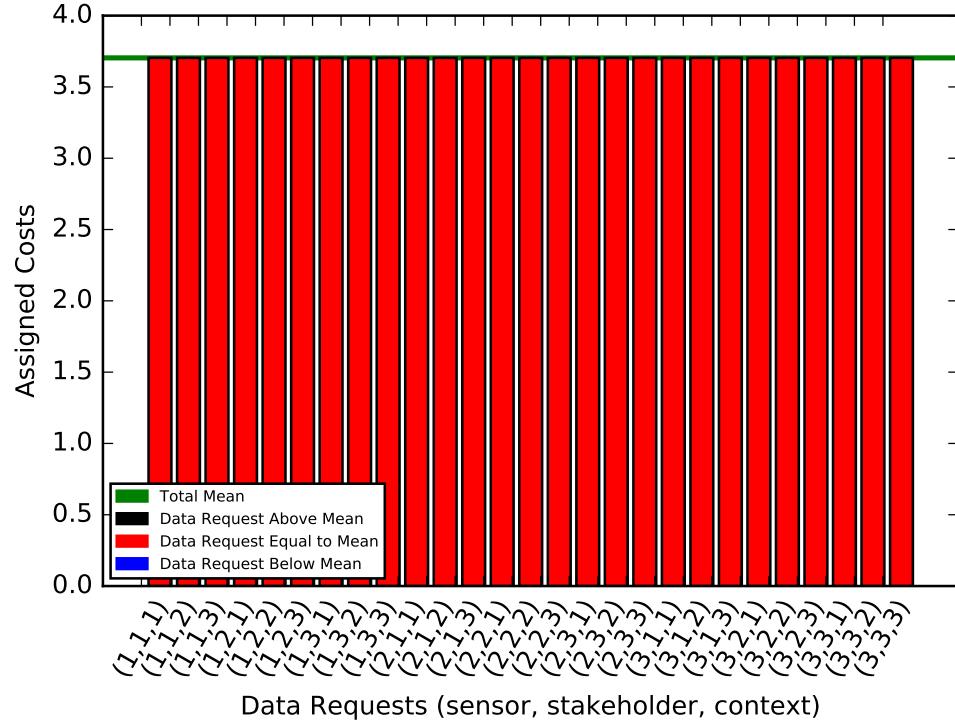


Figure 3.4: Values of the Cost Matrix

features in an equally intrusive way, then all the data requests will have the same weight and costs assigned.

Scenario 2

In this section we would like to test if data requests containing sub-features with higher intrusion levels are assigned higher weights and costs. Table 3.2 indicates the user input to scenario 2. As it can be seen, all features have equal categories, and all sub-features have the categories of 3 with an exception the Sensor's sub-features. The Sensors sub-features with identifiers 1,2 and 3 have respectively categories 1,3 and 5. This means that requests with Sensor's sub-feature 1 will be assigned a lesser weight in comparison to the other Sensors sub-features. Similarly, the data requests with Sensors sub-feature 2 will have a higher weightage assigned than Sensor's sub-feature 1 because of its higher category, but lesser than Sensor's sub-feature 3. Lastly, data requests with Sensor's sub-feature 3 will have a highest weight compared to the others, due to its category being 5. The weight and cost matrices can be seen in figures 3.5 and 3.6 respectively.

From the above inputs and graphs, we can conclude that the model assigns

3.3. Analysis of the Model

Table 3.2: Categorization for Scenario 2

Feature	Sub-Feature ID = 1	Sub-Feature ID = 2	Sub-Feature ID = 3
Sensors 3	Accelerometer 1	Noise 3	Location 5
Stakeholders 3	Corporation 3	Government 3	Educational Institution 3
Contexts 3	Navigation 3	Environment 3	Social Media 3

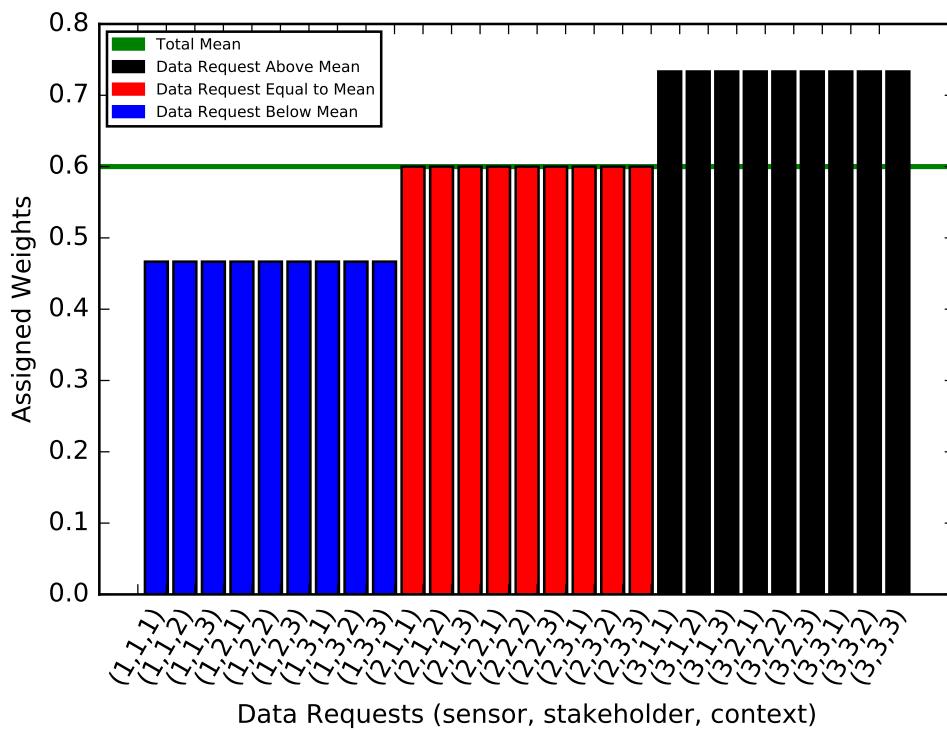


Figure 3.5: Values of the Weight Matrix

a higher weightage to data requests with sub-features that the user finds more intrusive compared to the others.

Scenario 3

The feature and sub-feature categories are both assigned different values in this section, to show how varying their values together affects the assignments of the weight and cost matrix. Table 3.3 is the user input to the scenario 3. All the features have different categories assigned from 3 to 5.

3. COMPUTATIONAL MODEL

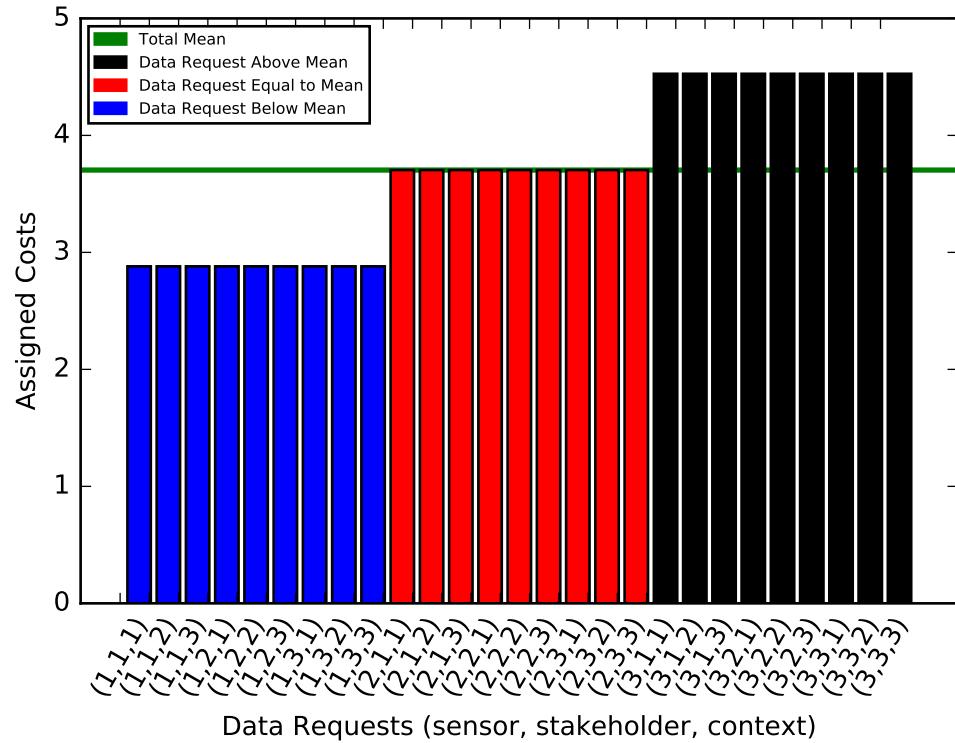


Figure 3.6: Values of the Cost Matrix

Additionally, the sub-feature 1 of each feature has a category of 5, higher than the other sub-features which are all categorized as 1. The weight and cost matrices generated for this scenario can be seen in figures 3.7 and ?? respectively.

Table 3.3: Categorization for Scenario 3

Feature	Sub-Feature ID = 1	Sub-Feature ID = 2	Sub-Feature ID = 3
Sensors	Accelerometer	Noise	Location
5	5	1	1
Stakeholders	Corporation	Government	Educational Institution
4	5	1	1
Contexts	Navigation	Environment	Social Media
3	5	1	1

As it is observed in both figures, the data request with the highest weight is the one with tuple (1,1,1). This tuple indicates that the data request involves

3.3. Analysis of the Model

all sub-feature 1 of each feature. This happens because all of the sub-feature 1 are assigned a category of 5. The feature sensors and its sub-feature 1 are categorized as 5, so all the data requests with tuple $(1, *, *)$, where * is all the other possible sub-features from other features, are all above average as seen in figures ??, irrespective of the categories of the other feature's sub-features. This shows that assigning a higher category to a feature can lead to higher data request costs.

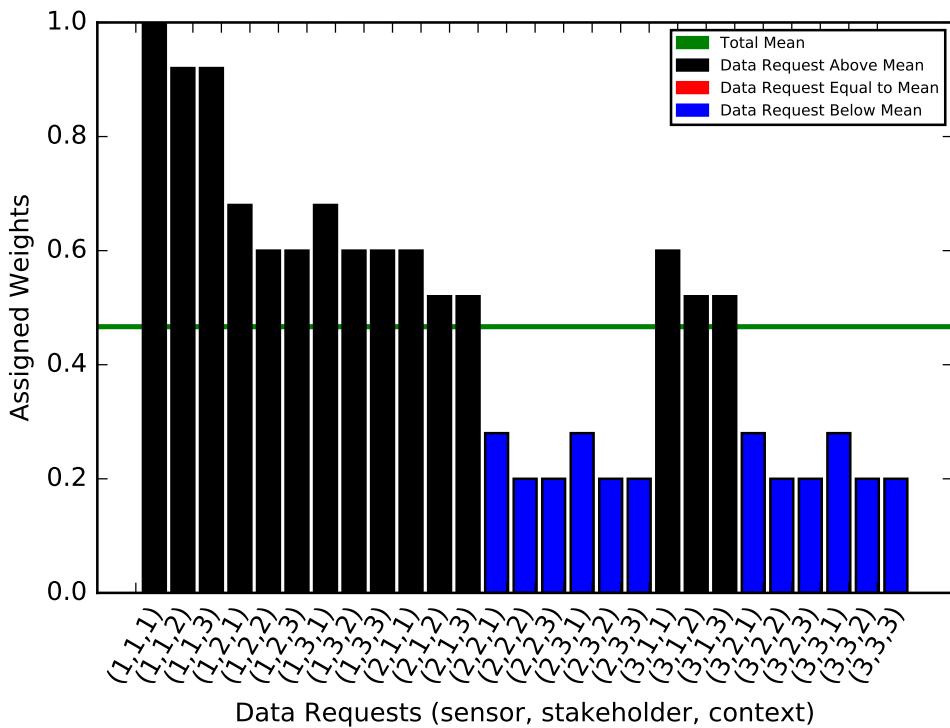


Figure 3.7: Values of the Weight Matrix

The green horizontal line in the graph indicates the mean value of the weights and costs. In general due to sub-features categorized as 5, those data requests receive a higher weight and cost. In some cases, the data requests still receive a lower weight such as tuple $(2,2,1)$, $(2,3,1)$, $(3,2,1)$ and $(3,3,1)$ even though Contexts sub-feature 1 has a category of 5. This is due to the fact that Sensors and Stakeholders feature have a higher category of 5 and 4 respectively than the contexts feature. Since their sub-features are assigned a lower privacy intrusion category than the context's sub-features, the weight of the data requests is lower. This shows that even though a sub-feature may be regarded as very intrusive, its weight increasing changing ability still depends on the category of the feature it belongs to.

3. COMPUTATIONAL MODEL

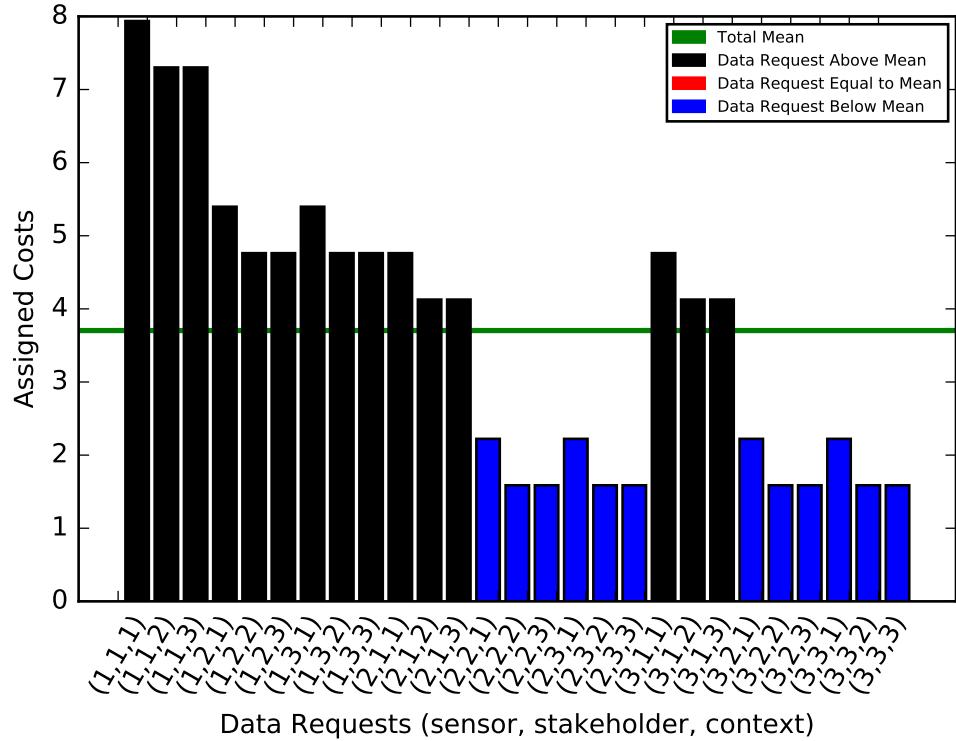


Figure 3.8: Values of the Cost Matrix

Additionally, it can be noted that data requests with at least two sub-features 1 are all above average. We can witness the property of the model, which puts more emphasis on the perception of the features than the sub-features themselves. As seen in the figure, all the features with higher intrusion categorizations have weights and costs that are well above average.

It can be concluded that the model assigns weights to data requests, by putting more emphasis on the feature's weights. A feature with high category has the ability to assign higher costs with a highly categorized sub-feature. It also has the ability to lower the weight of a data request with a sub-feature lowly categorized. Features with lower categories contribute lesser to the weight assignments, irrespective of their sub-feature categories.

Chapter 4

Experiment Methodology

In the previous chapter, the computational model has been explained in detail. This model has been implemented as a mobile application for the Android platform and can be used to collect real data from users. This application will help us collect information that can aid to see the influence of incentives on the data sharing decision. In this chapter we explain some of the work and decisions that are taken before and after the start of the experiment. It is then proceeded to explain how the experiment is carried out along with detailed instruction to the usage of the mobile application created.

4.1 Preparatory Phase

4.1.1 Pre-Survey

The pre-survey¹ is a survey created that runs before the deployment of the social experiment. This survey was made in order to study the perception of users on the three features to be studied which are explained in detail in section 4.2.2. Figure ??² depicts the features and their sub-features visually.

As it can be seen in the figure, there were a lot of sub-features to choose from each feature. Increasing the number of sub-features for each feature in the experiment in turn increases the number of data requests posed to the user. Additionally, we wanted to gain insight into the perception of users on the three features. Hence the survey was prepared to understand all of the above. Additionally, it can help us redesign some of the aspects of the experiment based on the ambiguities found and user feedback. The participants pool consist of both people who are aware and unaware of data

¹https://descil.eu.qualtrics.com/SE/?SID=SV_0xGS6kfmr8GtQd7

²Figure made by Athina Voulgari

4. EXPERIMENT METHODOLOGY

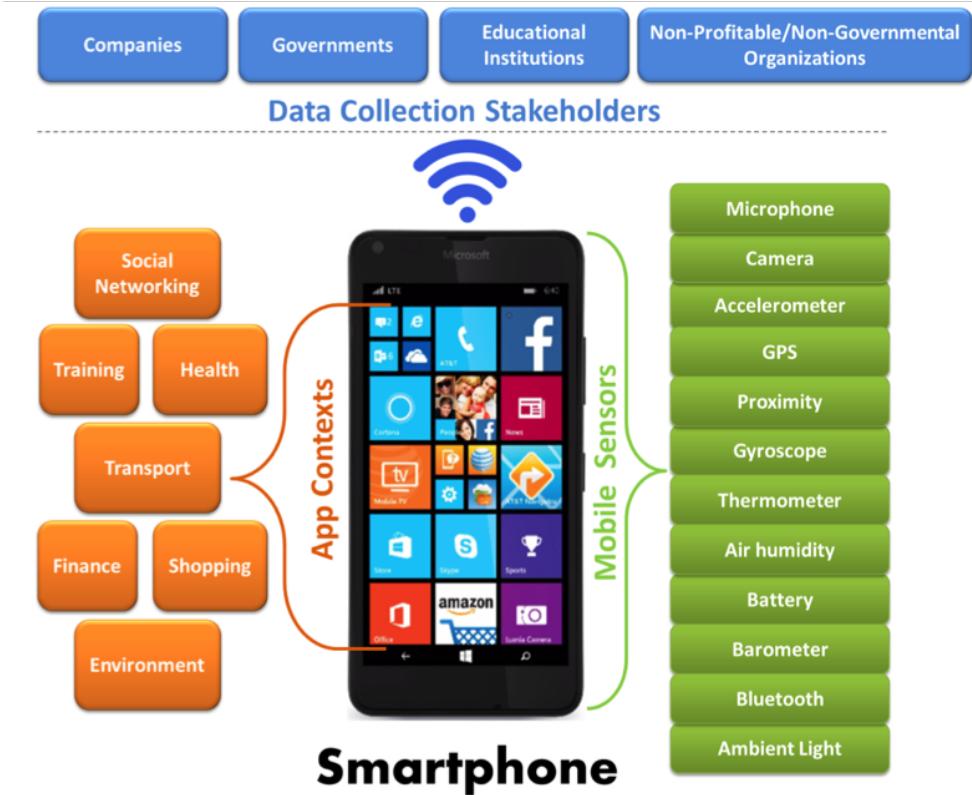


Figure 4.1: The Three Features Examined

privacy and sensors. Participants were not paid for filling out the survey. Till now, 199 entries have been recorded.

4.1.2 Sub-Features

Figures 4.2, 4.3 and 4.4, each show the average intrusion level of each possible sub-feature for the sensors, stakeholders and contexts. For the experiment, it was decided to choose for each feature two non-intrusive and two intrusive sub-features each. The minimum privacy intrusion level is one which indicates this sub-feature to not be intrusive, and the maximum is five which means that the sub-feature is very privacy intrusive.

For the sensors feature, it can be observed that the sub-features GPS and microphone are found to be have an intrusion of 4.2 and 3.8, which means users find these sensors on average very intrusive. On the other hand, sub-features light and accelerometer are found to be lower in intrusion with values of 2.2 and 2.3, which means that users find these sensors non-intrusive in general. The average of all sensors intrusion values is 2.8 as indicated by the blue line.

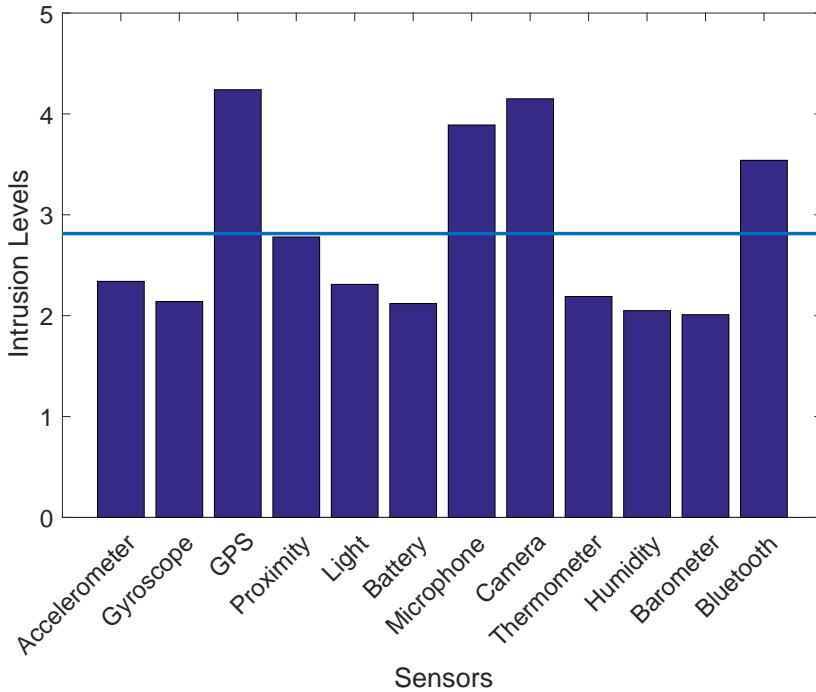


Figure 4.2: Average Intrusion of Sensors Sub-Features

Similarly, looking at the stakeholder feature graph 4.3, it is seen that sub-features corporation and government are found to be intrusive by the users with levels 3.8 and 3.6. On the other hand, sub-features educational institution and non governmental organization are found to be relatively less intrusive by the users with values of 3.2 and 2.95. For intrusion levels of contexts feature in graph 4.4, it is observed that sub-features social-networking and health are found to be intrusive by the users with values 3.8 and 3.6. Sub-features environment and transportation are regarded as less intrusive by user with values of 2.9 and 3.3. The above mentioned sub-features for every feature have been chosen for the experiment.

4.1.3 Privacy Options

Each data request is accompanied with privacy options ranging from 1 to 5 as explained in section 3.2.6. Option 1 indicates that the users would like to share their raw data without any sort of summarization or reduction in information. Option number 5 indicates that the users would not like to share their data for this data request. The options in between have linearly scaled summarization levels assigned to them ranging from least privacy (1) to most privacy (5). For more information on the summarization levels for

4. EXPERIMENT METHODOLOGY

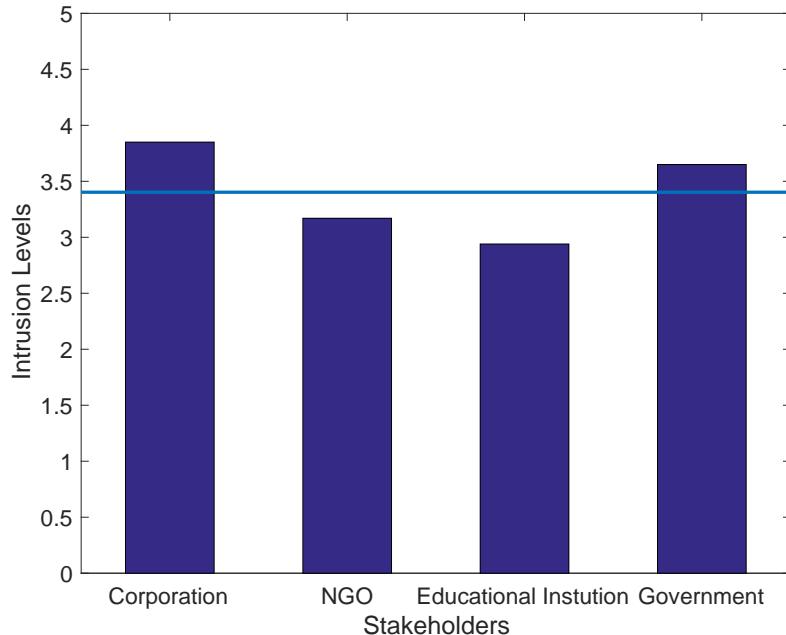


Figure 4.3: Average Intrusion of Stakeholders Sub-Features

each option please refer to section 3.2.8.

4.1.4 Question Structure

A data request is when a stakeholder asks users mobile sensor data for a particular context or purpose. Each data request to the user is posed in the form of a question with the following template :

"Please choose the amount of X sensor type data shared with Y stakeholder for use in a Z context app"

where Sensors X can be :

1. Accelerometer
2. Noise
3. Location
4. Light

where Stakeholders Y can be:

1. Corporation
2. Educational Institution

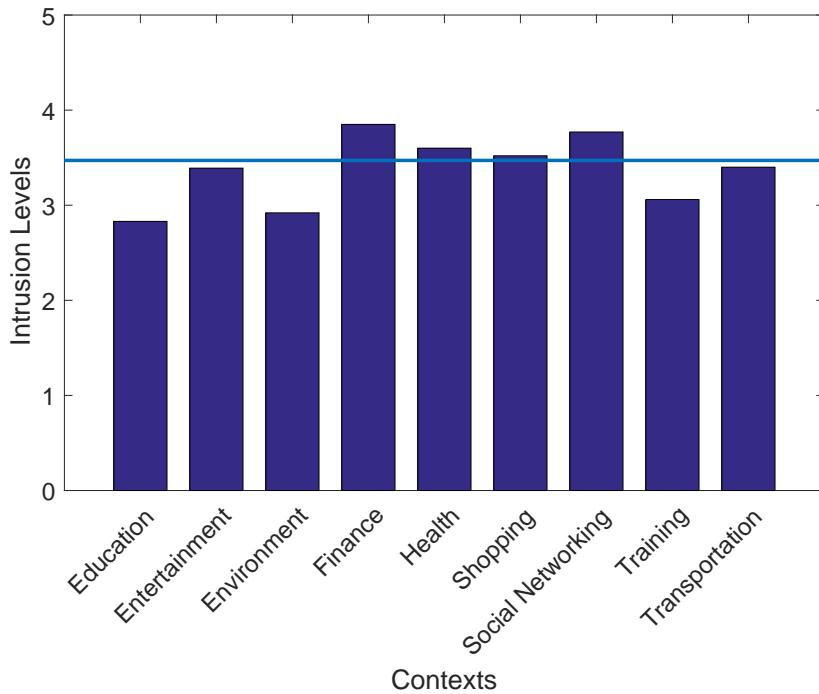


Figure 4.4: Average Intrusion of Contexts Sub-Features

3. Non Governmental Organization

4. Government

and where Contexts Z can be:

1. Environment
2. Health/Fitness
3. Navigation
4. Social Networking

In total this makes 64 data requests to the user. From now on, we will refer to mobile sensor data as just data.

4.1.5 Budget and Experiment Duration

The experiment is set to run for a total of two days, excluding the time taken for the entry phase and exit phase. The budget set for the core phase of the experiment is $b = 35$ Chf and is excluding the cost of participation in the entry and exit phase. Participants are paid 10 Chf for coming to the Entry Phase, and 15 Chf for participating in it. Similarly for the Exit Phase,

4. EXPERIMENT METHODOLOGY

participants are given 10 Chf for showing up, and 5 Chf for participating in it. Out of the budget B , $\frac{1}{7}$ is given away for the participation of the users in the core phase.

4.2 Entry Phase

The entry phase denotes the first day of the experiment. Users are asked to install the application from the PlayStore.

4.2.1 Collecting General User Information

As the figure 4.5 shows, the users are asked to answer some personal non-intrusive questions. The following is asked from the users:

1. Gender
2. Employment Status
3. Education Level
4. Year of birth
5. Country where user has lived most of his life
6. How many time a day do you check your Mobile phone per day.
7. Kind of applications the user has in the mobile phone.

The users may go back and re-answer the questions, but once the submit button is pressed on the screen ??, the data is sent to the server and hence cannot be changed. Users cannot navigate to the next pages without filling out all the questions.

4.2.2 Categorization of Features

As described in chapter 3, the users are asked to categorize the features sensors, stakeholders and contexts. As shown in figure 4.6a, each of the features are indicated followed by a drop down list of privacy options ranging from "*very low privacy intrusion*" to "*very privacy high privacy intrusion*". The option "*very low privacy intrusion*" means that the feature does not affect the users mobile sensor data sharing decision, whereas "*very privacy high privacy intrusion*" refers to a feature that very much affects the sharing of mobile sensor data.

Users need to click on the drop down menu to choose one of the privacy intrusion options. All the options are compulsory, and no default option is provided. Users cannot navigate to the next page without filling out all of the questions.

Figure 4.5: User Information Screens

4.2.3 Categorization of Sub-Features

For each of the features categorized in the previous sub-section, their sub-features need to be categorized in a similar fashion. Once again, the privacy options range from "*very low privacy intrusion*" to "*very high privacy intrusion*" like in section 4.2.2 . The users are first presented with the categorization of Sensors sub-features as shown in figure 4.7a.

Below each sensor is a drop down menu where the user can choose how much each of the sensors would affect the mobile sensor data sharing. Once all the sensors have been associated with a privacy intrusion level, the user can click the green submit button and is directed to the next page where the sub-features of stakeholders need to be in turn categorized in a similar fashion. This is depicted in figure 4.7b.

Each stakeholder type has a drop down menu each where the user can once again classify how much each of them affect data sharing. Once the user has finished entering the privacy intrusion level for stakeholders sub-features,

4. EXPERIMENT METHODOLOGY

Categorize Features

How intrusive are the following features of information sharing:

Sensors	very high privacy intrusion
Data Collectors	medium privacy intrusion
Context / Purpose	very low privacy intrusion

SUBMIT

GetUserInfo

Which types of apps do you usually have on your smartphone?

<input type="checkbox"/>	Educational
<input type="checkbox"/>	Entertainment
<input checked="" type="checkbox"/>	Finance
<input checked="" type="checkbox"/>	Games
<input checked="" type="checkbox"/>	Health&Fitness
<input type="checkbox"/>	Transport&Navigation
<input checked="" type="checkbox"/>	Music&Audio
<input type="checkbox"/>	News
<input type="checkbox"/>	Productivity
<input checked="" type="checkbox"/>	Shopping
<input checked="" type="checkbox"/>	Social Networking

SUBMIT

(a) Categorizing Features
(b) User Information Screen 3

Figure 4.6: Categorization and User Information Screens

the user can click the green submit button and is directed to the next page.

On this page, the users are asked to categorize how much each of the contexts sub-features affect mobile sensor data sharing. This is depicted in figure 4.8. Each context has a drop down menu below, where the user can rate each context. Once this has been done the user can click on the green submit button. The user will be redirected to the next page only if all the drop down boxes have been filled out. All questions are compulsory there is no default choice.

4.2.4 Answering Questions with No Incentives

After the categorization questions are answered and user answers are recorded, users will be presented with 64 questions. Each of these questions is a mobile sensor data request to the users. Users can choose from the available five privacy options mentioned in section 4.1.3. The options are indicated as a measure of how much data users can give, ranging from maximum data to least data. The higher the privacy of the option, the less information

4.3. Core Phase

The figure consists of two side-by-side screenshots of a mobile application. Both screens have a header bar at the top showing various icons and the time '23:30'. The left screen is titled 'Categorize Sensors' and asks 'How intrusive are the following sensors of information sharing?'. It lists five sensors: Accelerometer, Location, Light, Noise, and a fifth one whose name is partially visible. Each sensor has a corresponding orange rectangular button with a white border containing a text label: 'medium privacy intrusion' for Accelerometer, 'very high privacy intrusion' for Location, 'low privacy intrusion' for Light, and 'high privacy intrusion' for Noise. The fifth sensor's button contains the text '...'. At the bottom of this screen is a green rectangular button labeled 'SUBMIT'. The right screen is titled 'Categorize Stakeholders' and asks 'How intrusive are the following data collectors of information sharing?'. It lists four data collectors: Corporation, Non Governmental Organization, Government, and Education. Each collector has a corresponding orange rectangular button with a white border containing a text label: 'medium privacy intrusion' for Corporation, 'medium privacy intrusion' for Non Governmental Organization, 'high privacy intrusion' for Government, and 'very low privacy intrusion' for Education. At the bottom of this screen is a green rectangular button labeled 'SUBMIT'.

(a) Categorizing Sensors (b) Categorizing Stakeholders

Figure 4.7: Categorization of Sensors and Stakeholders Screen

about the sensor data is given away for that request and vice versa. Users can change the answers for a data request until the green submit button on top of the options that appears is clicked. The screen with the data request is shown in figure 4.9a.

After the users choose an option for the data request, a green submit button appears which is shown in figure 4.9b. Clicking on the submit button sends the response to the data request to the server and cannot be changed. At this stage, no indications of credit gained or privacy improvements are indicated.

Once all the questions have been answered, the user goes to the core phase of the experiment, which starts at day number two. In the experiment, day number one is the entry phase, the core phase is day number two and three.

4.3 Core Phase

Once the entry phase is done, the user is presented with the screen shown in figure ???. The "i" button at the bottom right of the screen denoted by

4. EXPERIMENT METHODOLOGY



Figure 4.8: Categorization of Contexts Screen

the number 9 is clickable. This takes the users to the FairDataShare portal. Figure ?? shows the homepage of the portal. Users can then click on the data generator registration section of the website where they can signup with their:

1. Username
2. Password
3. Email
4. Unique Identifier

The unique identifier is located at the bottom of the application screen is an alphanumeric sequence denoted by number 8. If it is long pressed the user can select the identifier, then copy and paste it in the textbox asking for the unique identifier in the portal. Figure ?? shows what the registration page looks like. The users can use this website to see all the data collected from them for all the mobile sensors. More details about the FairDataShare portal refer to the section 4.5.

The user can login into the portal after a minimum of 24 hours after the start of the core phase to see the data that has been collected and shared with the

4.3. Core Phase

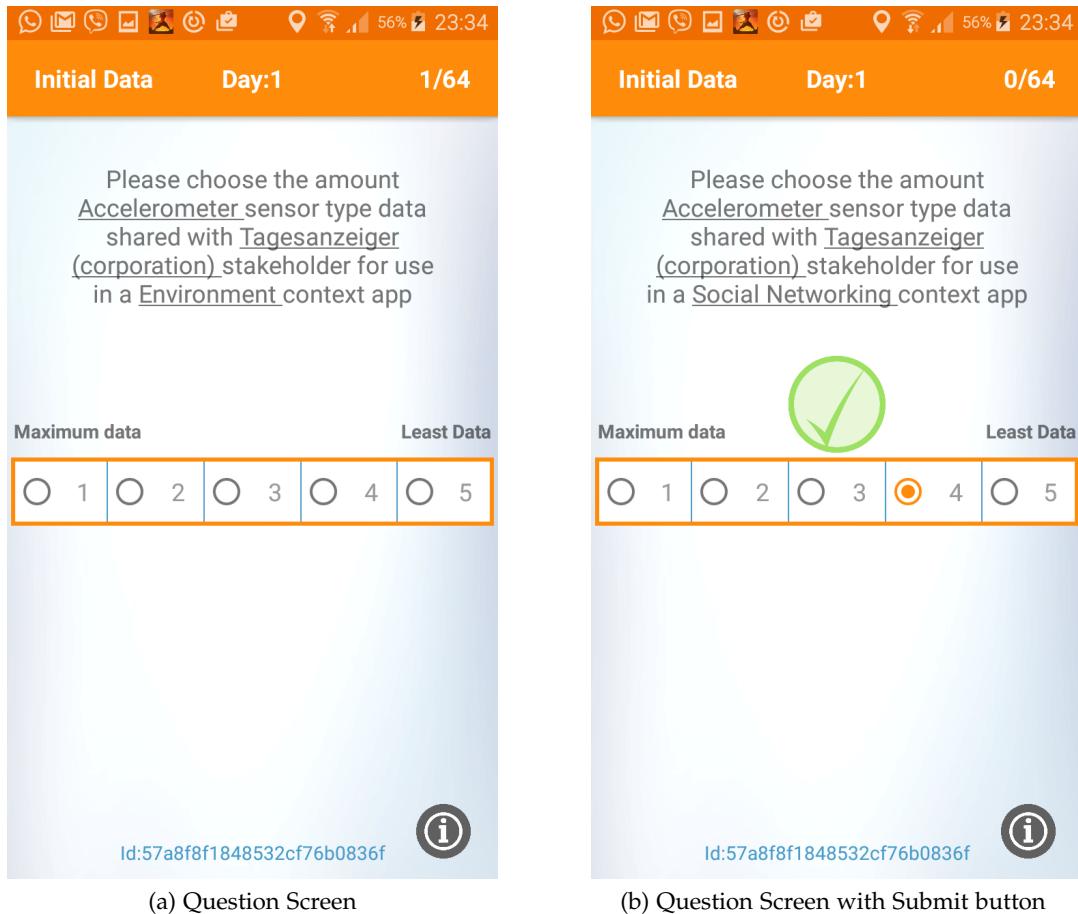


Figure 4.9: First Day Screen

stakeholders.

In the task-bar, the user can see the bidding day number and how many questions have been answered from the total available shown by numbers 1 and 2 in the figure ???. Day number one corresponds to the day where users answer questions with no incentives of any kind and was presented in the previous sub-section. The screen presented after the entry phase is over is what is called the "improvement screen". The button numbered 3 represents "improve privacy" and the button numbered 4 represents "improve credit" respectively. The items numbered 6 and 5 represent the privacy percentage and credit obtained by the user respectively. Privacy is measured in terms of the percentage of mobile sensor data not traded to the stakeholders. Credit is measured in terms of the currency Swiss Francs obtained for trading data to the stakeholders.

4. EXPERIMENT METHODOLOGY

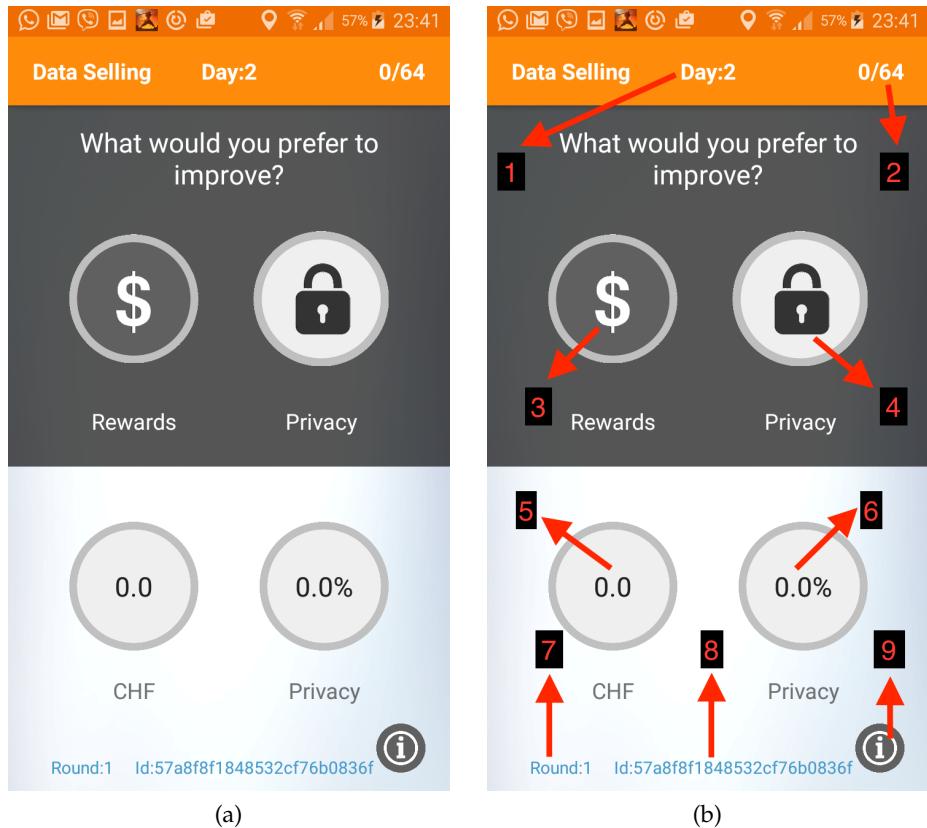


Figure 4.10: Improvement screen

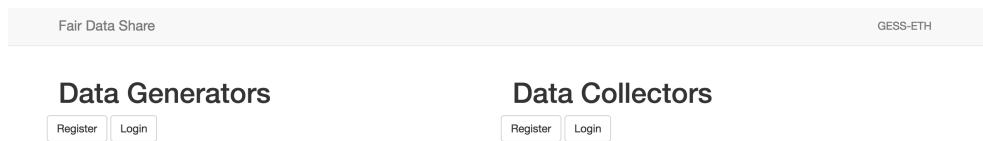


Figure 4.11: FairDataShare Homepage

4.3. Core Phase

Fair Data Share GESS-ETH

Data Generators

Check the code provided in the App

Register

Enter Username
Enter Password
Enter Email
Enter App Code

Register Cancel

Figure 4.12: Data Generators Registration Page

The item numbered 7 is the round number which indicates the number of times the user has answered all the data requests. The item numbered 2 is the number of questions the user has answered in the current round. Item number 1 indicates the experiment day number.

There are a total of 64 data requests, hence after all the 64 have been answered, the number of questions answered is reset and the number of rounds answered increases by one. This indicates all the data requests that have been answered and how many are left unanswered. Each question will have 5 options to choose from, ranging from maximum data sharing to least data sharing.

From the starting time of the core phase till 24 hours later marks one bidding day. Once 24 hours is over, another bidding day starts where the privacy and credit metrics are reset. The day number in the task bar is incremented by one. The user has to answer all the data requests again for this new bidding day. Previous responses to data requests are not carried over to the next day. If a data request is not answered, it is considered that the user does not want to trade mobile sensor data for that request. Additionally, each data request carries a participation fee, this is irrespective of the amount of mobile sensor data shared, by not participating in a data request the user foregoes this credit gain. The core phase goes on for a period of 48 hours.

4. EXPERIMENT METHODOLOGY

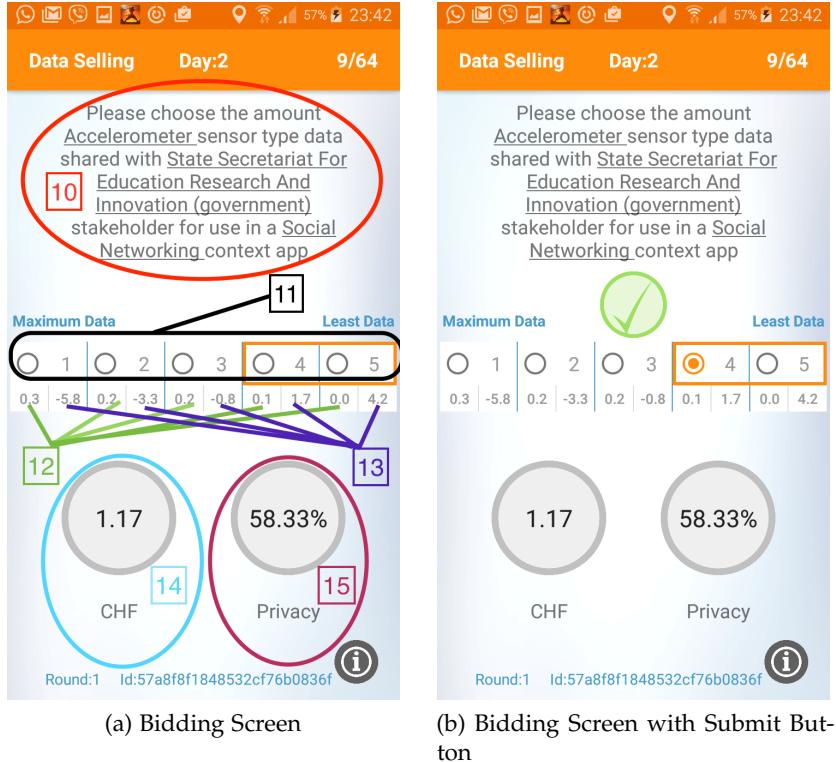


Figure 4.13: FairDataShare Portal

4.3.1 Improve Privacy or Credit

The improvement screen shown in figure 4.10 is where users can choose whether they would like to improve the privacy or the credit. The elements of this screen have been explained in the previous section 4.3. The improve credit button should be chosen if the user is interested in maximizing the amount of credit obtainable. This uses an algorithm that uses the previous user answers to put forth a data request that can increase the credit to the maximum explained in section 5.2.3. The credit improvement button is represented by the item number 5. Similarly, the improve privacy button is used to further improve the privacy that has been obtained. This puts forth a data request that can further increase the user privacy. It needs to be noted that the ultimate change in the privacy or credit metrics depends on the option chosen by the user for the data request. The privacy improvement button is represented by the number 6.

Scenario examples for each button is given in the next section after introducing the next screen in the application. For example, if a user chooses to improve the privacy, then clicks on improve privacy button and gets a data

request. The user still chooses option one with maximum data sharing (least privacy) for the data request, this may not improve his privacy but decrease it. This is because option 1 indicates that the user trades all the data for this request without filtering the sensor information. Trading all data gives the user more credit, but decreases the privacy metric.

Similarly, if a user chooses to improve the credit obtainable, the user clicks on the improve credit button and gets a data request. Then the user chooses the option five with least data sharing (maximum privacy) which indicates that no data is traded for this request. This response counters the initial desire to improve the credit obtainable. Trading no data increases one's privacy, but does not increase the credit to the maximum. Therefore, an actual improvement in the chosen metric depends on the chosen improvement button chosen and the choice of the appropriate option for that data request.

4.3.2 Answering Questions with Incentives

After choosing a metric to improve, a screen is presented as shown in figure 4.13a. This screen is called the "bidding screen". This screen is very similar to the screen 4.10 presented in the entry phase, except that the user is aware of the amount of privacy and credit obtained as indicated by items 14 and 15 respectively. Additionally, the user can see information about how the privacy and credit will increase or decrease for each privacy option of a data request. The items numbered 11 are the privacy options ranging from one to five.

The items numbered 12 are the improvement in privacy for each possible option of the current data request shown as item numbered 10. The items numbered 13 are the improvements in credit for each possible options of the current data request. Once the user decides on which options to choose according to how much data wants to be traded, the users can click on the radio option as explained in section 4.1.3 and then click again on the green submit button that pops up shown in 4.13b to confirm the answer. Once the green button has been clicked on, answers cannot be changed. The user has the possibility to go back to the improve screen from the bidding screen using the back button. Using the back button in the improve screen leads the user out of the application.

Additionally, for every question there is an orange recommendation box surrounding some options. This recommendation is highlighted by the number 16 in figure 4.14a. This gives an indication to the user as to which options can improve the privacy or the credit compared to the previous time the user has answered this data request. For example, if the user has previously answered option 4 to a data request and has clicked on improve credit, the system puts an orange box around options 1,2,3 and 4. Similarly, if the user clicked on improve privacy button, and the users previous answer was

4. EXPERIMENT METHODOLOGY

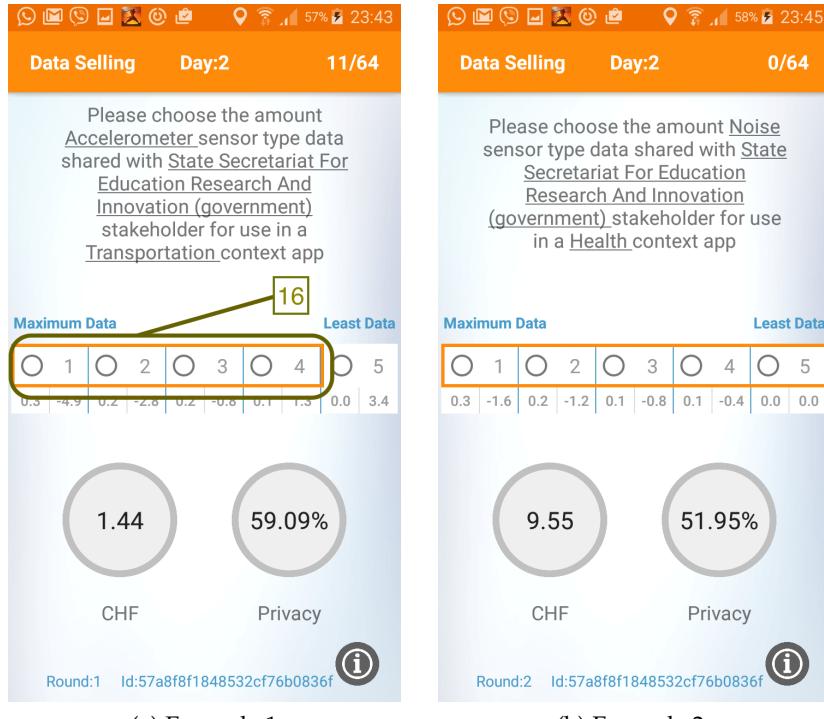


Figure 4.14: Recommendation Box

option 1, the system would recommend the options 1,2,3,4 and 5. Two examples of this are provided in figures 4.14.

It needs to be noted that the orange box does not necessarily provide an improvement of the particular metric chosen, it is meant to indicate improvements compared to the previous time the data request was answered to.

4.4 Exit Phase

After the end of the core phase, the participants are asked to fill up a survey based on their experience in the experiment. Some questions are about the rewards received, the privacy and credit metrics, design of the application, and how the experiment was conducted. The survey³ is linked to the user using the unique identifier assigned in the application. Once the survey is filled, the users receive their money for the entry phase, core phase and exit phase together, but only if they did not have their phones switched off throughout the experiment and participated in the core phase. This is done by checking the data collected on the server.

³https://descil.eu.qualtrics.com/SE/?SID=SV_3P0ySMqNe006v5j

4.5 FairDataShare Web Portal

The FairDataShare portal ⁴ is a website where users can view the data collected from them during the core phase of the experiment. Below is an explanation of how users and stakeholders can view mobile sensor data.

4.5.1 Data Generator's Portal

Once the users are registered which was explained in section 4.3, they can come back to the portal after a 24 hours period or later to view their mobile sensor data collected in the server. The data portal login page is shown in figure 4.15a. Since the users are already registered from the mobile phone in the entry phase, they can go to the portal from their computers and this time login instead of register. Users should enter their:

1. Username
2. Password

Once this is done, users will be redirected to the data collection page shown in figure 4.15b with the following options in the task-bar to choose from:

1. Accelerometer
2. Light
3. Noise
4. Location

Users can choose the sensor from the task-bar whose data they want to see by clicking on it. The data displayed includes the following columns :

1. Timestamp
2. Bidding day
3. Sensor Values

Figures 4.16a, 4.16b, 4.17a and 4.17b show examples of the data that can be seen for the location, light, accelerometer and noise sensor.

Users first register as data generators as indicated in the section 4.2.4.

4.5.2 Stakeholder's Portal

For a stakeholder to view data, they need to register in the portal shown in figure 4.11 by clicking register. Once that is done, the page in figure 4.18a is shown asking for the following details :

⁴<http://fair-data-share.inn.ac/>

4. EXPERIMENT METHODOLOGY

(a) Login Page

(b) Welcome Page

Figure 4.15: Entering the Portal

(a) Location Data

(b) Light Data

Figure 4.16: User Data

(a) Accelerometer Data

(b) Noise Data

Figure 4.17: User Data

1. Company Name
2. Email
3. Stakeholder Category
4. Company Website

The stakeholder category is the type the stakeholder comes under such as :

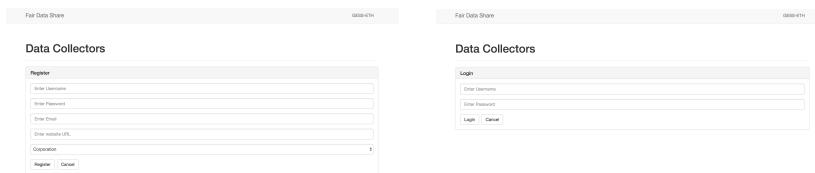
1. Corporation
2. Educational Institution
3. Government

4. Non-Governmental Organization

Once these details have been filled in, the stakeholder can click on the register button. Once registered, the stakeholder can login like shown in figure 4.18b. When access is granted the stakeholder is redirected to the page shown in figure 4.19. The stakeholder can choose from each of the available drop down lists :

1. A sensor
2. A context
3. An anonymous user
4. A bidding day number

Once this is entered, the stakeholder can see the data for that user with the privacy level decided by the anonymous user. If the stakeholder does not see any data, it means the user did not share data for this particular request. Stakeholders can view the sensor data in a similar fashion to users shown in figures 4.16 and 4.17. Data is available to the stakeholders 24 hours after the start of the core phase.



(a) Registration Page

(b) Login Page

Figure 4.18: Entering the Portal for Data Collectors

4. EXPERIMENT METHODOLOGY



Figure 4.19: Data Collectors Welcome Page

Chapter 5

Explanation of the Mobile Application

This chapter explains the details behind the making of the mobile application environment. First an overview is given, followed by a detailed explanation of the main components of the Android mobile application. This includes the architecture, database schemas and algorithms. Next, the server business logic and storage of the application is presented.

5.1 The Building Blocks

The following sections will explain integral parts of the server and client of the mobile application. A gist of the architecture is shown in the figure 5.1. As it can be seen, the mobile application represents the user participating in the experiment. As the experiment goes on, mobile sensor data and responses to the data requests which are collected are periodically sent to the Kinvey Data Store. The users can choose to login into the FairDataShare Portal from their computer or the mobile application. Once the user is authenticated, the user requests are sent from the FairDataShare server to the Kinvey Data Store¹. Kinvey in turn fetches the appropriate data and gives it back to the FairDataShare Server. This in turn structures the data so it can be easily readable, and pushes it to the user to see on the portal. The concept is similar for the Stakeholders, except they can only access the portal through the computer and not the mobile application.

5.2 The Mobile Application

The mobile application was developed for the Android platform with phones having API above level 17². Phones are assumed to have internet connectiv-

¹<https://kinvey.com>

²<https://developer.android.com/guide/topics/manifest/uses-sdk-element.html>

5. EXPLANATION OF THE MOBILE APPLICATION

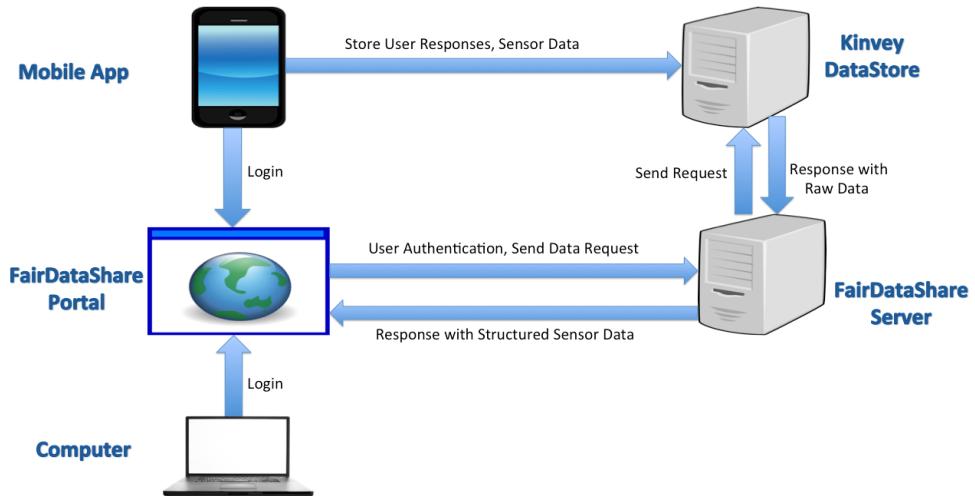


Figure 5.1: Conceptual Diagram of Mobile Application Architecture

ity and sufficient storage space of at least 100 Mb. Below is an explanation of some of the tasks that take place in the application.

5.2.1 Local Storage

The local storage is an integral part of the application. The database used is SQLite³ and is the default database for the Android environment. Small sized unrelated data pieces are stored in preference files (as key value pairs), whereas larger related data is stored in the database. The following paragraphs will explain each table present in this application followed with their function and schema. All tables explained here are pertaining to the user using the mobile application and not the server.

Figure 5.2a shows the **QUESTION_STORE**'s table schema. This table stores each possible data request with its features such as with its sensor **SENSOR**, stakeholder **STAKEHOLDER** and context **CONTEXT**. Each of these are represented by an integer, for example sensor 0 stands for accelerometer sensor. Each data request is accompanied by an unique question identifier **QID**, weight assigned **WEIGHT** and the cost assigned **COST**. This data is not sent to the server.

Figure 5.2b depicts the table **WHICH_ANSWERS**'s table schema. This stores the questions identifier **QID** of each data request that has been answered by the user for each round. This is helpful while fetching data requests, so as not to fetch the request twice in the same round. It makes sure that all questions

³<https://developer.android.com/reference/android/database/sqlite/package-summary.html>

5.2. The Mobile Application

QUESTION_STORE	WHICH_ANSWERED
<ul style="list-style-type: none"> 👉 Q_ID: INTEGER ▢ SENSOR: INTEGER ▢ STAKEHOLDER: INTEGER ▢ CONTEXT: INTEGER ▢ COST: REAL ▢ WEIGHT: REAL 	<ul style="list-style-type: none"> 👉 Q_ID: INTEGER

(a) Table Schema of QUESTION_STORE's
(b) Table Schema of WHICH_ANSWERED

Figure 5.2: Table Schemas

are answered before answering them for a second time. This data is not sent to the server.

STORE_ANSWERS	STORE_POINTS
<ul style="list-style-type: none"> 👉 Q_ID: INTEGER ▢ LEVEL: INTEGER ▢ DAY: INTEGER ▢ COST_OBT: REAL 	<ul style="list-style-type: none"> 👉 DAY: INTEGER ▢ PRI: REAL ▢ COST: REAL

(a) Table Schema of STORE_ANSWERS
(b) Table Schema of STORE_POINTS

Figure 5.3: Table Schemas

Figure 5.3a explains the schema of STORE_ANSWERS table. This table is used to store the data request identifier *QID* with the corresponding user responses *LEVEL*, along with the increase or decrease in credit obtained *COST_OBT*. The total cost can be calculated by adding all the costs in this table. Similarly, the total privacy can be calculated by averaging of all the user responses stored in this table. Only the most recent responses are stored in this table. The content of the table is not sent over to the server.

Figure 5.3b denotes the schema of STORE_POINTS table. This table is used to store the credit and privacy obtained for each bidding day. This information is sent to the server as soon one bidding day is over.

5. EXPLANATION OF THE MOBILE APPLICATION

USERRESPONSE_CACHE	
KEY:	INTEGER
UR:	VARBINARY(2000)
IS_SENT:	INTEGER

Figure 5.4: Table USERRESPONSE_CACHE Schema

Figure 5.4 depicts the USERRESPONSE_CACHE table's schema. This table stores a unique key *KEY* for each user response, followed by a flag *ISSENT*, which is 1 if the response is not sent to the server, and 0 if it is sent. The user response saved consists of the following entries :

1. User Identifier
2. Timestamp of the response
3. Sensor Identifier
4. Stakeholder Identifier
5. Context Identifier
6. Privacy Level response for this data request
7. Cost obtained for this data request
8. Current Total Privacy of the user
9. Current Total Credit of the user
10. Maximum Obtainable Credit for this data request in this round
11. Metric Chosen to Improve (Improve Privacy or Improve Credit)

All of the above fields are packed into the field *ur* shown in 5.4. The data in this table is sent to the server. Once the entry is sent to the server, the *ISSENT* field is changed to 0 and deleted locally. The unique keys *KEY* are useful for deleting sent entries. Figure 5.5 and 5.6 show the table schemas for data storage of the following sensors:

1. Accelerometer in the STORE_ACCELEROMETER table
2. Noise in the STORE_NOISE
3. Location in the STORE_LOCATION

5.2. The Mobile Application

4. Light in the STORE_LIGHT

The general schema for all the sensor tables is the following :

1. *KEY* - Uniquely identifies each sensor entry
2. *TIMESTAMP* - The time the sensor value was collected
3. *ISSENT* - Denotes whether the sensor entry has been sent to the server or not
4. The other columns are specific to each sensor and represent the actual sensor values collected

STORE_ACCELEROMETER	STORE_NOISE
<ul style="list-style-type: none">👉 KEY: INTEGER⌚ X: REAL⌚ Y: REAL⌚ Z: REAL⌚ TIMESTAMP: NUMERIC(15,0)⌚ IS_SENT: BOOLEAN	<ul style="list-style-type: none">👉 KEY: INTEGER⌚ RMS: REAL⌚ SPL: REAL⌚ BANDS: CHARACTER(20)⌚ TIMESTAMP: NUMERIC(15,0)⌚ IS_SENT: BOOLEAN

(a) Table Schema of STORE_ACCELEROMETER (b) Table Schema of STORE_NOISE

Figure 5.5: Table Schemas for Sensor Data

STORE_LOCATION	STORE_LIGHT
<ul style="list-style-type: none">👉 KEY: INTEGER⌚ LAT: REAL⌚ LONG: REAL⌚ TIMESTAMP: NUMERIC(15,0)⌚ IS_SENT: BOOLEAN	<ul style="list-style-type: none">👉 KEY: INTEGER⌚ X: REAL⌚ TIMESTAMP: NUMERIC(15,0)⌚ IS_SENT: BOOLEAN

(a) Table Schema of STORE_LOCATION (b) Table Schema of STORE_LIGHT

Figure 5.6: Table Schemas for Sensor Data

5. EXPLANATION OF THE MOBILE APPLICATION

5.2.2 Alarms and Notifications

Every bidding day where the user can answer data requests lasts for a period of 24 hours. After one bidding day is over, the system needs to be informed in a timely manner to perform some application critical functions. The function performed are explained in detail in section 5.2.2. To inform the system of such an event Android provides the functionality in the form of alarms.

Alarms can be set to go off just once or in a repeated fashion to trigger tasks. Unfortunately, the alarms provided by Android are not exact for some versions ⁴, in the sense that they are triggered around that time set but not exactly to optimize the battery, and can be delayed upto 24 hours. Hence, it is decided to set the repeating alarms manually.

The first time the application opens the alarm is set to ring in exactly 24 hours, but things change when the phone is switched off. One of the conditions of the experiment is not to have the phone switched off at any time. Nevertheless, it is taken into account the scenario where the phone is kept switched off for a period of time. There are various things that can happen:

1. The phone is rebooted.
2. The phone is switched off, during this time an alarm is missed.
3. The phone is switched off for a period greater than 24 hours. One or more alarms can be missed.

Once the phone is switched off, all alarms are erased from memory ⁵. Alarms do not execute when the phone is switched off. Hence, when the phone switches on, BootReceiver service of the application is triggered with pseudocode shown in 1. This checks whether an alarm has been missed, if it has been missed 200 seconds is given for the phone to stabilize after boot before triggering tasks. Otherwise, a new alarm is set using the pseudocode shown in 2. To set an alarm we need the time difference between now and when the alarm should ring. After that is calculated, the alarm is set.

Going to the Next Data Sharing Day

Once the alarm rings, it marks the end of a bidding day. Once a bidding day ends a number of tasks need to be executed and for this the NextDayService is triggered, which is described in pseudocode shown in 5. To start with the privacy and credit is sent to the sent to the server and stored locally in the STORE_POINTS table. *Privacy* which is the total privacy obtained, *Credit* is the total credit obtained, *Round* which is the number of time the user answered

⁴<https://developer.android.com/training/scheduling/alarms.html>

⁵<https://developer.android.com/reference/android/app/AlarmManager.html>

Algorithm 1 BootService Algorithm

```

1: procedure BOOTSERVICE
2:   now  $\leftarrow$  current timestamp
3:   i  $\leftarrow$  timestamp of last triggered alarm
4:   if now  $-$  i  $<$  86400 then
5:     Call SetAlarmLater()
6:   else
7:     Set alarm in 200 seconds

```

Algorithm 2 Alarm Algorithm

```

1: procedure SETALARMLATER
2:   now  $\leftarrow$  current timestamp
3:   i  $\leftarrow$  timestamp of last triggered alarm
4:   latertime  $\leftarrow$  i + 86400
5:   latergap  $\leftarrow$  latertime  $-$  now
6:   Set Alarm in latergap seconds

```

all the questions and *CurrentQuestion* which is the current question the user is answering is all reset to zero. The *Day* corresponds to the current day number is incremented by one to denote the next bidding day.

Algorithm 3 NextDayService Algorithm

```

1: procedure NEXTDAYSERVICE
2:   Store Privacy, Credit, Day in STOREPOINTS
3:   Send Privacy, Credit, Day to Server
4:   Privacy, Credit, Round, CurrentQuestion  $\leftarrow$  0
5:   Day  $\leftarrow$  Day + 1
6:   Store current time
7:   Call Summarization()
8:   if Day  $>$  End then
9:     End experiment
10:  else
11:    Update user interface elements

```

The current time of executing the alarm is saved in case the phone is rebooted or switched off. After that, the sensor data which is saved locally needs to be summarized, the corresponding method is called and is explained in pseudocode shown in 4. Finally it needs to be checked if the experiment is over or not and update the user interface accordingly. This means either the various metrics on the improvement and bidding screens (which ever is currently active) are updated, or the end of experiment screen

5. EXPLANATION OF THE MOBILE APPLICATION

is shown.

5.2.3 Fetching Data Requests

A data requests need to be fetched from the database in two scenarios :

1. After a question has been answered in the first bidding day (entry phase)
2. After the privacy or credit improvement button has been clicked (core phase)

In the first bidding day, once a data request has been answered the next one is fetched sequentially from the database. This just requires knowing the current data request number and fetching the next data request from table QUESTION_STORE. For the other bidding days, fetching of the data requests depends on the improvement button chosen. According to the choice, the following is done:

1. **Improve Privacy** - Obtain data request from table STORE_ANSWERS where user has answered with lowest privacy
2. **Improve Credit** - Obtain data request from table STORE_ANSWERS where user has answered with highest privacy

In addition to sending the data request to the user interface, we need to show how choosing each option of the data request will affect the total privacy and total credit metrics. To do this for the total cost, we output the computation *last – possible*, where *last* stands for the credit obtained the last time the data request was answered. *possible* stands for the maximum amount of credit that can be obtained for this option (each data request has five privacy options 3.2.6). The possible total cost changes are shown under the options. For more detail on how credits are split among options in a data request refer 4.1.3.

Every option of a data request has an associated percentage of data that is given away as described in 3.2.6. According to the percentage of data given away, the total privacy is calculated for each possible option. The difference between the current privacy and each possible total privacy is calculated and indicated under each option. This gives an indication to the user as to what each option will do to the metrics.

5.2.4 Recording User Choices

The figure 5.4 describes the table USERRESPONSE_CACHE. Each time a user enters a response to a data request, all the fields mentioned in section 5.2.1 are recorded and stored in a class object. This object is transformed into a byte array so as to be stored easily in the table as is without transformation.

5.2. The Mobile Application

When the JobNetworkService described in 5.2.6 is called, the class object is sent as it is to the server after converting it back to an object.

5.2.5 Sensor Data Collection and Summarization

Sensor data is collected from the following sensors :

1. Accelerometer sensor
2. Noise sensor
3. Location sensor
4. Light sensor

A sensor service is triggered when the application is installed and is stopped when the experiment is over. This collects data from every sensor every 30 seconds and stores it in the appropriate tables mentioned in section 5.2.1. At the end of a bidding day, sensor data needs to be summarized according to the wishes of the user. This starts by first finding out the lowest privacy level for each sensor. Privacy levels range from one to five, that is from the lowest to highest privacy levels. Using this level summarization is done as shown in pseudocode 4. Every privacy level corresponds to an action:

1. 1- All data is sent to the server
2. 2- Send 75% of the data
3. 3- Send 50% of the data
4. 4- Send 25% of the data
5. 5- Do not send any data

Initially all the sensor data has a field *ISSENT* with value of zero. Data that should be sent to the server is set with *ISSENT* = 1, and all others that have value *ISSENT* = 0 are ignored.

5.2.6 Server Synchronization

User responses and sensor data need to be sent to the server. This is done periodically every 5000 seconds in order to free up space on the phone whenever the internet is available. It is triggered first when the application is started for the first time. Data is fetched from the tables in the database. Data with fields marked as *ISSENT* = 1 is data that is ready and that has not been sent yet to the server. Such data is sent, and when an acknowledgement is received from the server, this data is deleted from the table.

5. EXPLANATION OF THE MOBILE APPLICATION

Algorithm 4 Summarization Algorithm

```
1: procedure SUMMARIZATION
2:   for each sensor do
3:     Fetch sensor data from sensor table
4:     level  $\leftarrow$  Fetch user privacy level
5:     if level  $\leftarrow$  1 then
6:       Set all ISSENT  $\leftarrow$  1
7:     else if level  $\leftarrow$  2 then
8:       for 3 out of every 4 records do
9:         ISSENT  $\leftarrow$  1
10:    else if level  $\leftarrow$  3 then
11:      for 1 out of every 2 records do
12:        ISSENT  $\leftarrow$  1
13:    else if level  $\leftarrow$  4 then
14:      for 1 out of every 4 records do
15:        ISSENT  $\leftarrow$  1
16:    Delete all entries with ISSENT  $\leftarrow$  0
17:    Update Database
```

Algorithm 5 JobNetworkService Algorithm

```
1: procedure NETWORKSERVICE
2:   Fecth data from USERRESPONSECACHE
3:   for each record do
4:     if ISSENT == 1 then
5:       Send record to Server
6:       if SUCCESS then
7:         Delete record
8:   for each sensor do
9:     Fecth data from sensor table
10:    for each record do
11:      if ISSENT == 1 then
12:        Send record to Server
13:        if SUCCESS then
14:          Delete record
```

5.3 The Server

5.3.1 Kinvey Data Storage

Kinvey⁶ is a mobile backend as a service which provides a platform for mobile phones to link applications to a backend cloud storage⁷. For the purpose of this application the backend has been used to store data and for some business logic implementations in javascript.

Security

All communications from the application to the server is encrypted using TLS/SSL encryption⁸ to communicate with the backend service. This is automatically provided and done by the Kinvey SDK.

Collection Store

Locally, all information is stored in SQLite which is a relational database. The database used in Kinvey is MongoDB so instead we have collections on the server. When the user starts the application, general personal information is entered as explained in 5.2.1. This data is stored in the collection UserInformation with the schema shown in the screen shots 5.7 and 5.8.

birth_year	check_mobile_frequency	country	education	education_background	education_level	employment_status	entertainment	finance
1994	3	"France"	0	0	3	6	0	1
1924	3	"Arménie"	0	0	4	4	0	0
1923	3	"Armenia"	0	0	4	2	0	0
1922	3	"Aruba"	0	0	3	2	0	0
1992	1	"France"	0	0	4	6	0	0
1991	1	"France"	0	0	5	6	0	0
1924	3	"Argentin..."	0	0	3	2	0	0
1921	3	"Andorre"	0	0	2	1	0	0
1923	3	"Argentin..."	0	0	2	2	0	0
1922	3	"Antigua..."	0	0	2	1	0	0
1922	3	"Anguilla"	0	0	2	2	0	0
1923	3	"Anguilla"	0	0	2	1	0	0
1923	3	"Angola"	0	0	2	2	0	0
1926	2	"Angola"	0	0	4	6	0	0
1924	3	"Andorre"	0	0	3	5	0	0
1924	3	"Angola"	0	0	4	6	0	0
1920	5	"Fiji"	0	0	1	3	0	0

Figure 5.7: Screenshot of Collection UserInformation Part 1

⁶<http://kinvey.com/>

⁷https://en.wikipedia.org/wiki/Mobile_backend_as_a_service

⁸Kinvey white paper : KINVEY CLOUD SERVICE: SECURITY OVERVIEW 2014

5. EXPLANATION OF THE MOBILE APPLICATION

gender	health	medical	mobile_sensor_privacy	music	user_id	navigation	news	productivity	shopping	social_network
2	1	0	3	1	"57a8f8f1848532cf7...	0	0	0	1	1
2	0	0	3	0	"579a148f352257bc0...	0	0	0	0	0
2	0	0	3	1	"57975541e813f9973...	0	0	0	0	0
2	0	0	2	0	"57975159890927b61...	0	0	0	0	0
2	0	0	3	1	"57935b55a67b0ba32...	1	0	0	0	0
2	1	0	3	1	"579357faa67b0ba32...	1	0	0	0	1
2	1	0	3	0	"579357faa67b0ba32...	0	0	0	0	0
1	0	0	3	0	"57931cad866a46bd5...	1	0	0	0	0
2	1	0	3	0	"57930d55837af5db6...	0	0	0	0	0
1	1	0	3	0	"5792471c493006891...	0	0	0	0	0
2	1	0	3	0	"57923fbfc3d7cee30...	0	0	0	0	0
2	1	0	3	0	"57923946c3d7cee30...	0	0	0	0	0
2	0	0	3	0	"5792373d3692318e3...	1	0	0	0	0
2	0	0	4	1	"57922e1afb5591741...	0	0	0	0	0
2	1	0	3	0	"57922a329bb316492...	0	1	0	0	0
2	0	0	3	1	"579224ffba4636590...	0	0	0	0	0
2	0	0	1	0	"5791d082bb71b5202...	0	0	0	0	0
1	1	0	4	1	"578e91e778f251171...	1	1	0	0	1

Figure 5.8: Screenshot of Collection UserInformation Part 2

Once this is done, users have to categorize the various Features, Sensors, Stakeholders and then the various Contexts. This information is sent to the server in collections named Features, Sensors, Stakeholders and Contexts. Schema is shown in 5.9, 5.10, 5.11 and 5.12 respectively.

user_id	context	data_collector	sensor
"57a8f8f1848532cf7...	1	3	5
"579a148f352257bc0...	2	3	3
"57975541e813f9973...	4	5	1
"57975159890927b61...	4	3	1
"57935b55a67b0ba32...	1	3	5
"579357faa67b0ba32...	1	3	5
"57931cad866a46bd5...	3	2	2
"57930d55837af5db6...	4	1	2
"5792471c493006891...	4	2	2
"57923fbfc3d7cee30...	4	1	2
"57923946c3d7cee30...	3	1	2
"5792373d3692318e3...	4	1	2
"57922e1afb5591741...	4	2	2
"57922a329bb316492...	3	2	4
"579224ffba4636590...	4	2	2
"578e91e778f251171...	3	5	4
"578e28687d1cdd1b6...	4	3	2

Figure 5.9: Screenshot of Collection Features

All the data stored locally on the mobile phone which is sent by the JobNet-

5.3. The Server

user_id	acc	gps	light	noise
"57a8f8f1848532cf76b0836f"	3	5	2	4
"57a8f8f1848532cf76b0836f"	3	5	2	4
"579a148f352257bc0612c70b"	2	2	4	3
"57975541e813f99735dd0598"	2	4	3	3
"57975159890927b613d0f49f"	1	3	5	3
"57935b55a67b0ba32f81eeac"	3	5	1	5
"579357faa67b0ba32f81e724"	3	5	1	5
"57931cad866a46bd554a1896"	2	2	3	4
"57930d55837af5db6734a438"	2	1	2	4
"5792471c493006891e4e0c2b"	2	2	3	4
"57923fbfc3d7cee306b50957"	2	3	3	4
"57923946c3d7cee306b4fb44"	2	4	2	4
"5792373d3692318e3ba9e929"	2	1	2	4
"57922e1afb55917415770d44"	2	3	4	4
"57922a329bb316492f70242b"	2	1	2	4
"579224ffba46365901e66e78"	2	2	2	4
"578e91e778f2511711cfb9f5"	1	5	1	3

Figure 5.10: Screenshot of Collection Sensors

user_id	corp	edu	gov	ngo
"57a8f8f1848532cf76b0836f"	3	1	4	3
"579a148f352257bc0612c70b"	3	3	4	2
"57975541e813f99735dd0598"	1	4	1	3
"57975159890927b613d0f49f"	2	4	4	5
"57935b55a67b0ba32f81eeac"	5	3	5	5
"579357faa67b0ba32f81e724"	5	3	5	5
"57931cad866a46bd554a1896"	2	4	4	1
"57930d55837af5db6734a438"	3	3	1	2
"5792471c493006891e4e0c2b"	2	4	1	1
"57923fbfc3d7cee306b50957"	3	4	3	2
"57923946c3d7cee306b4fb44"	3	3	3	3
"5792373d3692318e3ba9e929"	3	3	2	2
"5792373d3692318e3ba9e929"	3	3	2	2
"57922e1afb55917415770d44"	4	3	2	1
"57922a329bb316492f70242b"	2	4	2	1
"579224ffba46365901e66e78"	2	4	2	1
"578e91e778f2511711cfb9f5"	5	2	5	4

Figure 5.11: Screenshot of Collection Stakeholders

workService explained in section 5.2.6 is received by Kinvey. User responses are stored in the collection UserResponse shown in 5.13 and 5.14.

The sensor data sent by the JobNetworkService is stored in collections named after the sensors themselves. The schema of the tables is shown in figures 5.15, 5.16, 5.18 and 5.17.

5. EXPLANATION OF THE MOBILE APPLICATION

user_id	environment	health	social_networking	transportation
"57a8f8f1848532cf76b0836f"	3	3	5	3
"579a148f352257bc0612c70b"	4	2	2	2
"57975541e813f99735dd0598"	2	2	4	2
"57975159890927b613d0f49f"	3	1	5	5
"57975159890927b613d0f49f"	3	1	5	5
"57975159890927b613d0f49f"	3	1	5	5
"57935b55a67b0ba32f81eeac"	1	1	1	1
"579357faa67b0ba32f81e724"	3	1	5	3
"57931cad866a46bd554a1896"	3	3	3	5
"57930d5837af5db6734a438"	3	2	1	3
"5792471c493006891e0c2b"	1	2	1	3
"57923fbfc3d7cee306b50957"	3	3	2	5
"57923946c3d7cee306b4fb44"	3	2	3	3
"5792373d3692318e3ba9e929"	2	3	1	3
"5792373d3692318e3ba9e929"	2	3	1	3
"5792373d3692318e3ba9e929"	2	3	1	3
"5792373d3692318e3ba9e929"	3	3	2	4

Figure 5.12: Screenshot of Collection Contexts

contexts	credit	credit_can_be	credit_gain	credit_question	data_collectors	timestamp	day_no
0	6.510009765625	0	-0.07690429687499999	0.17944335937499994	0	"2016-07-25 11:51:57.234"	3
1	6.436767578125	0.30029296874999994	0	0.35034179687499994	2	"2016-07-25 10:19:00.938"	3
0	6.436767578125	0.30029296874999994	0	0.35034179687499994	2	"2016-07-24 14:53:15.508"	3
2	6.436767578125	0.30029296874999994	0	0.35034179687499994	1	"2016-07-24 14:53:04.694"	3
3	6.436767578125	0.30029296874999994	0	0.35034179687499994	1	"2016-07-24 14:53:13.376"	3
1	6.436767578125	0.30029296874999994	0	0.35034179687499994	1	"2016-07-24 14:52:55.396"	3
0	6.436767578125	0.30029296874999994	0	0.35034179687499994	1	"2016-07-24 14:52:41.651"	3
2	6.436767578125	0.30029296874999994	0	0.35034179687499994	0	"2016-07-24 14:52:37.05"	3
3	6.436767578125	0.30029296874999994	0	0.35034179687499994	0	"2016-07-24 14:52:39.684"	3
1	6.25732421875	0.17944335937499994	0.17944335937499994	0.17944335937499994	2	"2016-07-24 14:23:08.071"	3
3	6.436767578125	0.17944335937499994	0.17944335937499994	0.17944335937499994	2	"2016-07-24 14:23:09.255"	3
1	6.436767578125	0.30029296874999994	0	0.35034179687499994	0	"2016-07-24 14:52:34.802"	3
0	6.436767578125	0.30029296874999994	0	0.35034179687499994	0	"2016-07-24 14:52:20.038"	3
0	6.077880859375	0.17944335937499994	0.17944335937499994	0.17944335937499994	2	"2016-07-24 14:23:07.005"	3
1	5.5395078125	0.17944335937499994	0.17944335937499994	0.17944335937499994	1	"2016-07-24 14:23:02.056"	3
3	5.8984375	0.17944335937499994	0.17944335937499994	0.17944335937499994	1	"2016-07-24 14:23:04.974"	3
2	5.718994140625	0.17944335937499994	0.17944335937499994	0.17944335937499994	1	"2016-07-24 14:23:03.856"	3

Figure 5.13: Screenshot of Collection UserResponse Part 1

To keep track of all the existing users in the experiment, the collection Users stores all unique user identification strings of participants. The schema is shown in 5.19.

Finally, the collection Score shown in 5.20 stores the total privacy, total credit obtained by the user for each bidding day.

5.3. The Server

Figure 5.14: Screenshot of Collection UserResponse Part 2

day_no	user_id	lat	long	summarization	timestamp
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186498206
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186468203
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186408196
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186438200
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186378192
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186288183
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186348189
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186258180
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186228177
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186318186
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186198171
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186168166
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186138163
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186108160
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186078154
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186018145
3	"578e91e778f2511711cfb9f5"	47.419864654541016	8.502890586853027	1	1469186048150

Figure 5.15: Screenshot of Collection Location

Bussiness Logic

?? Most of the business logic used for the FairDataShare portal is present in Kinvey. There are two main scripts stored in Kinvey:

- ## 1. Script to find the privacy preference

5. EXPLANATION OF THE MOBILE APPLICATION

day_no	summarization	timestamp	user_id	x	y	z
3	1	1469186493918	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.15141487121582
3	1	1469186463719	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.15141487121582
3	1	1469186373308	"578e91e778f251171...	0.0191536135971546...	-0.124498486518859...	10.180145263671875
3	1	1469186433709	"578e91e778f251171...	0.0287304203957319...	-0.134075298905372...	10.15141487121582
3	1	1469186343108	"578e91e778f251171...	0.0287304203957319...	-0.134075298905372...	10.15141487121582
3	1	1469186403508	"578e91e778f251171...	0.0191536135971546...	-0.134075298905372...	10.15141487121582
3	1	1469186282709	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.15141487121582
3	1	1469186222308	"578e91e778f251171...	0.0287304203957319...	-0.105344876646995...	10.20887565612793
3	1	1469186252509	"578e91e778f251171...	0.0287304203957319...	-0.143652096390724...	10.160991668701172
3	1	1469186312909	"578e91e778f251171...	0.0287304203957319...	-0.143652096390724...	10.15141487121582
3	1	1469186131909	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.160991668701172
3	1	1469186192307	"578e91e778f251171...	0.0287304203957319...	-0.134075298905372...	10.160991668701172
3	1	1469186162108	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.160991668701172
3	1	1469186101709	"578e91e778f251171...	0.0287304203957319...	-0.143652096390724...	10.15141487121582
3	1	1469186011488	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.132261276245117
3	1	1469185981478	"578e91e778f251171...	0.0287304203957319...	-0.143652096390724...	10.160991668701172
3	1	1469186071708	"578e91e778f251171...	0.0191536135971546...	-0.143652096390724...	10.15141487121582

Figure 5.16: Screenshot of Collection Accelerometer

bands	user_id	day_no	rms	spl	summarization	timestamp
"0.0,1.9080862E-5,...	"578e91e778f2511711cfb9f5"	3	107.31494140625	62.65440368652344	1	1469186468205
"0.0,1.5665331E-5,...	"578e91e778f2511711cfb9f5"	3	88.882568359375	61.01753234863281	1	1469186498214
"0.0,1.952634E-5,...	"578e91e778f2511711cfb9f5"	3	100.1298828125	62.05247497558594	1	1469186438208
"0.0,2.1573624E-5,...	"578e91e778f2511711cfb9f5"	3	47.744140625	55.61960220336914	1	1469186408210
"0.0,2.0647063E-5,...	"578e91e778f2511711cfb9f5"	3	73.727294921875	59.39376449584961	1	1469186378210
"0.0,2.1016425E-5,...	"578e91e778f2511711cfb9f5"	3	71.015380859375	59.0682487487793	1	1469186348202
"0.0,2.062715E-5,...	"578e91e778f2511711cfb9f5"	3	96.7724609375	61.7562370300293	1	1469186258193
"0.0,2.1374477E-5,...	"578e91e778f2511711cfb9f5"	3	67.723876953125	58.656036376953125	1	1469186168186
"0.0,1.7658736E-5,...	"578e91e778f2511711cfb9f5"	3	106.538818359375	62.59135818481445	1	1469186198197
"0.0,1.8755187E-5,...	"578e91e778f2511711cfb9f5"	3	40.89111328125	54.27377700805664	1	1469186318192
"0.0,1.429928E-5,...	"578e91e778f2511711cfb9f5"	3	21.7080078125	48.77359771728156	1	1469186078184
"0.0,1.9742341E-5,...	"578e91e778f2511711cfb9f5"	3	85.58349609375	60.68899917602539	1	1469186138189
"0.0,2.134739E-5,...	"578e91e778f2511711cfb9f5"	3	78.437744140625	59.93170166015625	1	1469186108185
"0.0,2.007436E-5,...	"578e91e778f2511711cfb9f5"	3	82.403076171875	60.360069274902344	1	1469186048176
"0.0,2.0650641E-5,...	"578e91e778f2511711cfb9f5"	3	76.4638671875	59.71023333740234	1	1469185958181
"0.0,1.4806586E-5,...	"578e91e778f2511711cfb9f5"	3	21.33837890625	48.624427795410156	1	1469186018178
"0.0,2.0179677E-5,...	"578e91e778f2511711cfb9f5"	3	101.021484375	62.12947463989258	1	1469185988177

Figure 5.17: Screenshot of Collection Noise

2. Script for summarization

The stakeholders make a request for data on the FairDataShare portal giving the following details:

1. Bidding day number
2. Anonymous user
3. Sensor

5.3. The Server

day_no	summarization	timestamp	user_id	x
3	3	1469447239362	"57935b55a67b0ba32...	47
3	3	1469447109437	"57935b55a67b0ba32...	54
3	3	1469447319071	"57935b55a67b0ba32...	39
3	3	1469446323286	"57935b55a67b0ba32...	109
3	3	1469446998112	"57935b55a67b0ba32...	180
3	3	1469446812605	"57935b55a67b0ba32...	165
3	3	1469446228120	"57935b55a67b0ba32...	83
3	3	1469445977205	"57935b55a67b0ba32...	96
3	3	1469445805362	"57935b55a67b0ba32...	156
3	3	1469445621109	"57935b55a67b0ba32...	157
3	3	1469445343373	"57935b55a67b0ba32...	136
3	3	1469445541953	"57935b55a67b0ba32...	143
3	3	1469445255903	"57935b55a67b0ba32...	150
3	3	1469444855996	"57935b55a67b0ba32...	127
3	3	1469444963549	"57935b55a67b0ba32...	127
3	3	1469445171668	"57935b55a67b0ba32...	100

Figure 5.18: Screenshot of Collection Light

user_id
"57a8f8f1848532cf76b0836f"
"579a148f352257bc0612c70b"
"57975541e813f99735dd0598"
"57975159890927b613d0f49f"
"57935b55a67b0ba32f81eeac"
"579357faa67b0ba32f81e724"
"57931cad866a46bd554a1896"
"57930d55837af5db6734a438"
"5792471c493006891e4e0c2b"
"57923fbfc3d7cee306b50957"
"57923946c3d7cee306b4fb44"
"5792373d3692318e3ba9e929"
"5792373d3692318e3ba9e929"
"57922e1afb55917415770d44"
"57922a329bb316492f70242b"
"579224ffba46365901e66e78"
"578e91e778f2511711cfb9f5"

Figure 5.19: Screenshot of Collection Users

4. Context

Given this input plus the category of the stakeholder (which is known from the registration), we look into the UserResponse Collection trying to find the most recent record that fits this criteria and extract the privacy level.

5. EXPLANATION OF THE MOBILE APPLICATION

timestamp	user_id	credit	day_no	privacy
"2016-08-08 23:35:20.788"	"57a8f8f1848532cf76b0836f"	6.310096153846153	1	74.609375
"2016-07-28 16:23:21.687"	"579a148f352257bc0612c70b"	9.030898876404493	1	56.25
"2016-07-26 14:22:12.236"	"57975541e813f99735dd0598"	9.01900773195876	1	57.03125
"2016-07-26 14:09:18.367"	"57975159890927b613d0f49f"	8.850940265486726	1	58.203125
"2016-07-26 14:00:00.565"	"57935b55a67b0ba32f81eeac"	0	4	0
"2016-07-25 13:59:00.084"	"57935b55a67b0ba32f81eeac"	6.510009765625	3	63.28125
"2016-07-24 13:59:31.599"	"57935b55a67b0ba32f81eeac"	8.8885498046875	2	50.390625
"2016-07-23 13:59:31.382"	"57935b55a67b0ba32f81eeac"	5.90576171875	1	67.1875
"2016-07-23 09:31:59.106"	"57931cad866a46bd554a1896"	1.2409156976744176	1	55
"2016-07-23 08:36:40.373"	"57930d55837af5db6734a438"	8.348214285714286	1	60.9375
"2016-07-22 18:18:33.16"	"5792471c493006891e4e0c2b"	1.2428977272727268	1	52.5
"2016-07-22 17:48:06.151"	"57923fbfc3d7cee306b50957"	1.253551136363636	1	55
"2016-07-22 17:19:40.016"	"57923946c3d7cee306b4fbb44"	9.385190217391308	1	53.125
"2016-07-22 16:32:46.545"	"57922e1afb55917415770d44"	7.999999999999998	1	63.28125
"2016-07-22 16:22:07.321"	"57922a329bb316492f70242b"	8.459821428571429	1	60.9375
"2016-07-22 13:22:00.152"	"578e91e778f2511711cfb9f5"	14.11313657407408	3	18.359375

Figure 5.20: Screenshot of Collection Score

Once we know the privacy level, summarization can be done. Data has been taken from the user with a certain summarization, and if the summarization level is lower than the privacy level extracted, further summarization needs to be done. The pseudocode is shown in 6.

Algorithm 6 Server Summarization Algorithm

```

1: procedure SUMMARIZATION
2:   data  $\leftarrow$  sensor data from collection
3:   if summarizationlevel == privacylevel then
4:     Return data
5:   else
6:     skip  $\leftarrow$  summarizationlevel - privacylevel + 1
7:     for every skipnumber records out of 4 do
8:       Delete record from data
9:   Return data to portal

```

5.3.2 FairDataShare Web Portal

The FairDataShare portal makes use of a server at ETH Zurich other Kinvey to safely store the usernames, passwords of the users and the stakeholders in a collection. The database technology used is MongoDB. The language used to interact with Kinvey is Express.js, which is based on Node.js. Most of the data portal business logic is on Kinvey as described in section ???. The

5.3. The Server

webpage was constructed using simple Html and css. All screenshots of the portal including detailed information is provided in chapter 4.

Chapter 6

Pre-Survey and Experiment Findings

The following chapter will give an overview of the data obtained from the survey, which was conducted before running the experiment. Later, an overview of the data obtained from the experiment is explained along with feedback received from the participants. This chapter puts forth all that was learnt from the above mentioned.

6.1 Overview of the Pre-Survey Data

The survey has 199 participants. After filtering out spurious and half-filled entries 189 entries are used for the data analysis. In the following paragraphs we will presenting the information obtained in the survey.

Out of the total participants 63.64% are male and 36.36% are female. The mean birth year was found to be 1985. The demographics of the participants is explained in the table 6.1. On the education level 2.53% have not completed high school, 9.60% have completed high school, 5.05% have gone to some college, 28.79% have obtained their bachelors degree, 39.90% have gotten their masters degrees and 14.14% have obtained their PHDs. About the employment of the participants 51.52% are full time employees, 6.06% are part time employed, 6.06% are unemployed and looking for work, 1.52% are unemployed and not looking for work, 0.51% are retired and 41.92% are students. None of the participants are disabled.

Figure 6.1 depicts the various kinds of applications the population has on their mobile phones. As it can be seen, the most popular applications are social networking, transportation and music applications. Plot 6.2 shows the frequency of mobile usage among the population. It is observed that the majority of the people use their phones 36-70 times a day.

In figure 6.3, the bars with the color dark blue represents for each level the percentage of participants. The scale of the levels are from one to five, one

6. PRE-SURVEY AND EXPERIMENT FINDINGS

indicates that people do not care to five indicating that people care a lot about mobile sensor data. As observed, most people have answered level 3, meaning they care to a medium level and 77.28% care to a level of 3 and above.

In the same figure, the color light blue shows the amount of contribution that sensors in general have in the data sharing decision. Level one indicates that it contributes none, level 5 means that it contributes a lot. It can be observed that most people find that sensors contribute to a level of 4. Similarly, the color green shows the amount of contribution stakeholders in general have in the data sharing decision. It can be observed that most people find that it contributed to a level of 4. Lastly, the color yellow depicts the level of contribution contexts have in the data sharing decision. Most people feel that it contributes to a level of 5.

From the above, it is understood that in general, more than the sensor, the more important thing in a data request is for what purpose the data is being collected followed by the entity to whom the data is traded with. The privacy intrusion levels of the participants for the individual different sensors, stakeholders and contexts has been explained in chapter 4.

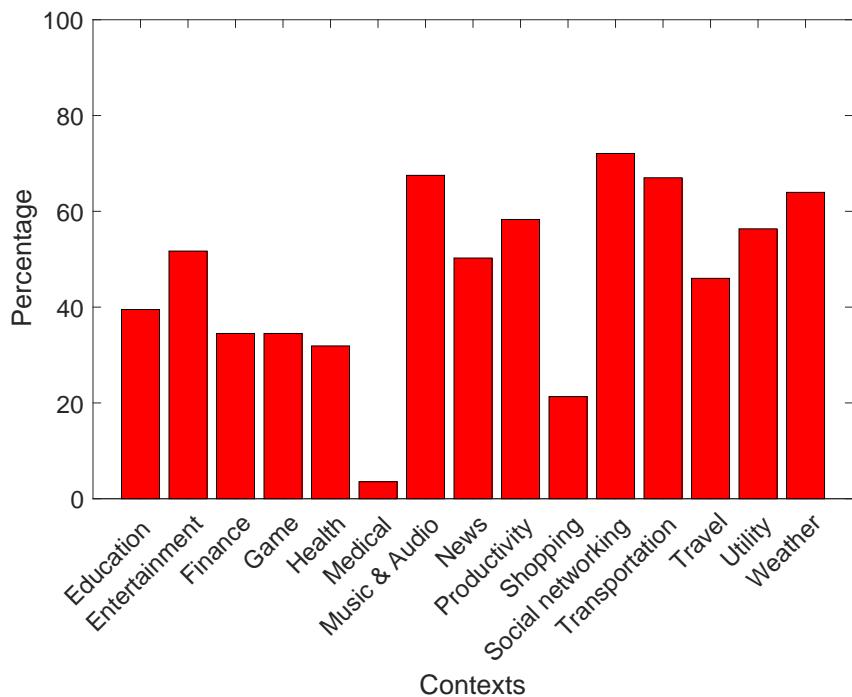


Figure 6.1: Applications in the Mobile Phone

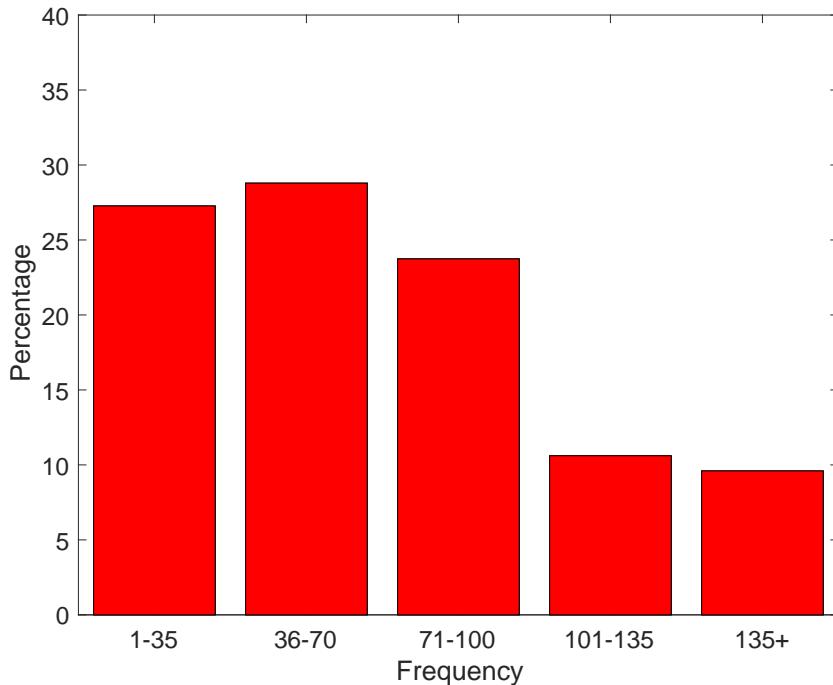


Figure 6.2: Frequency of Mobile Phone Usage

6.2 Pre-Survey Methodology and Findings

All the results presented below were performed on the data by performing the following changes to the data:

1. Rows with empty fields were removed
2. Rows with spurious data were removed
3. Data was scaled or normalized when necessary

Other than the above, the data was not manipulated. Outliers were not excluded either.

Perception of Individual Sensor Grouped on the Intrusion of Sensors in General

We try to examine here if the perception of intrusion of Sensors can affect the way a person views the individual sensors themselves. In other words, we try to examine if there is a significant difference in perception of each sensor depending on the perception of the Sensors as a whole. For this, we grouped the survey data based on the responses to question 10. Since there

6. PRE-SURVEY AND EXPERIMENT FINDINGS

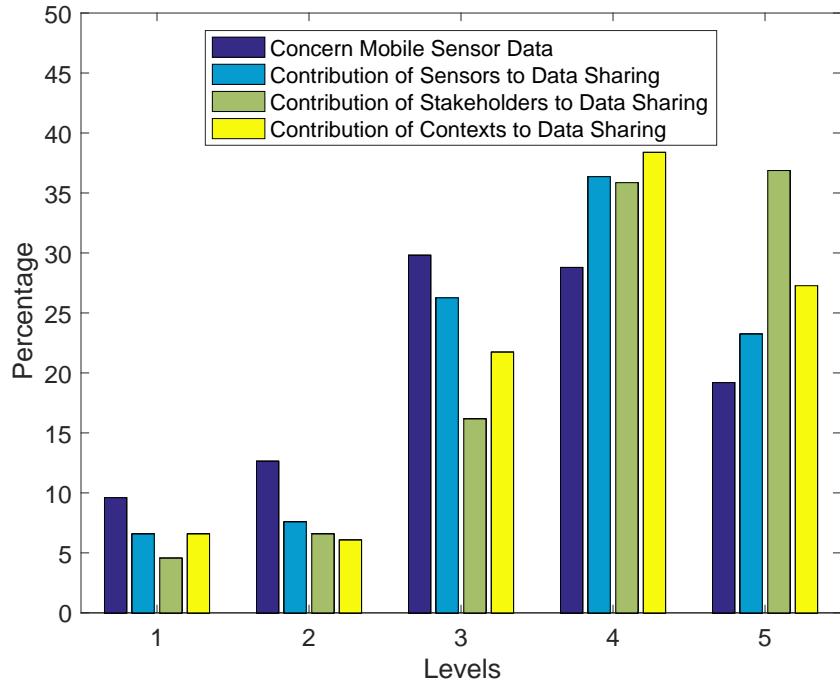


Figure 6.3: Graph depicting the concern of Mobile Sensor Data and the contribution of various features to the data sharing decision

are 5 possible responses to this question, this makes 5 individual groups from 1 to 5.

We now have 5 groups who view sensors in a different light each and their perception of each of the individual sensors can be compared. Before going into the comparison, we try to understand the properties of the data. To perform a one-way ANOVA test or a t-test, the data needs to be:

1. Normally distributed
2. Homoscedastic
3. Ordinal or continuous

Since the data is discrete and follows the Likert Scale with options from 1 to 5, it gives skewed normal distribution. Additionally, the variances of values within the groups formed are not similar. One-way ANOVA test is quite robust to heteroscedacity, as long as the maximum variance among all groups is less than four times the group with the lowest variance. The scale used to collect data is in the ordinal form. Accounting for all the violations, we instead opt for a non-parametric tests such as the Kruskal-Wallis H test and the Dunn's test which only assume the following :

6.2. Pre-Survey Methodology and Findings

Table 6.1: Demographics of Population

Country	Percentage
United States of America	1.01%
United Arab Emirates	0.51%
The former Yugoslav Republic of Macedonia	0.51%
Syrian Arab Republic	0.51%
Switzerland	20.71%
Spain	1.01%
Slovakia	0.51%
Serbia	5.05%
Russian Federation	0.51%
Netherlands	1.52%
Italy	2.02%
Iran	1.01%
India	14.65%
Hungary	0.51%
Greece	29.29%
Germany	10.61%
France	1.52%
Czech Republic	1.01%
Costa Rica	0.51%
China	0.51%
Columbia	0.51%
Canada	0.51%
Bolivia	0.51%
Brazil	1.52%
Bahrain	0.51%
Argentina	0.51%
Austria	2.02%

1. Groups are independant from one another
2. All observations are independant
3. The dependant variables should be in the ordinal scale or continuous

The above tests do not make any assumptions about the distribution of the data and are robust to heteroscedastic data.

Group 1 to 5 have 13, 14, 50, 71 and 42 people each respectively. For in depth analysis of the composition of each group in terms of employment, education, gender and birth year please refer to tables 6.2,6.5,6.3,6.4. Figures 6.4a and 6.4b depict the mean and variances for each of the individual groups. We start by performing the non parametric Kruskal-Wallis test on

6. PRE-SURVEY AND EXPERIMENT FINDINGS

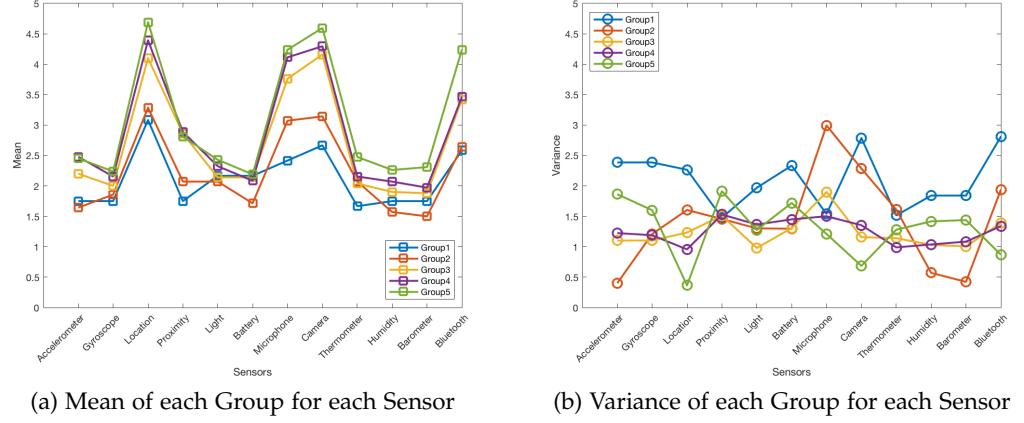


Figure 6.4: Table Schemas

each Sensor. The value of alpha assumed here is 0.05. The null hypothesis states that all the groups perceive the sensors in the similar way. This means they come from the same distribution. The alternative hypothesis is that the groups perceive each sensor in a significantly different way. The table 6.6 depicts the p-values obtained from the test.

Table 6.2: Employment Classification of Groups

Occupation	1	2	3	4	5
Employed full time	4.90%	6.86%	26.47%	38.24%	23.53%
Employed part time	8.33%	16.67%	33.33%	16.67%	25.00%
Unemployed and looking for work	8.33%	16.67%	16.67%	25.00%	33.33%
Unemployed and not looking for work	0.00%	0.00%	0.00%	66.67%	33.33%
Retired	0.00%	0.00%	100.00%	0.00%	0.00%
Student	7.23%	4.82%	26.51%	39.76%	21.69%
Disabled	0.00%	0.00%	0.00%	0.00%	0.00%

Table 6.3: Gender Classification of Groups

Gender	1	2	3	4	5
Female	5.56%	4.17%	33.33%	30.56%	26.39%
Male	7.14%	9.52%	22.22%	39.68%	21.43%

On these sensors, we proceed with a post hoc test by performing a pairwise Dunn's test to examine if there is an actual significant difference between

6.2. Pre-Survey Methodology and Findings

Table 6.4: Average Birth Year of Groups

1	2	3	4	5
1989	1979	1986	1986	1983

Table 6.5: Education Classification of Groups

Education	1	2	3	4	5
Less than high school	20.00%	20.00%	20.00%	40.00%	0.00%
High school	10.53%	0.00%	52.63%	31.58%	5.26%
Some college	10.00%	20.00%	30.00%	20.00%	20.00%
Bachelors degree	7.02%	10.53%	24.56%	42.11%	15.79%
Masters degree	3.80%	5.06%	24.05%	35.44%	31.65%
PhD degree	7.14%	7.14%	17.86%	35.71%	32.14%

Table 6.6: Kuskal-Wallis Test

Sensor	p-value
Accelerometer	0.0151
Gyroscope	0.2959
Location	1.0664e-05
Proximity	0.0147
Light	0.6933
Battery	0.6950
Microphone	3.0070e-04
Camera	2.1191e-05
Thermometer	0.0693
Air Humidity	0.1292
Barometer	0.0949
Bluetooth	3.4877e-05

the groups and if so between which groups. The sensors with p values with less than 0.05 are examined in more detail and the p-values are presented in table 6.7. The table shows the results for each pairwise test done, with the p-values adjusted using the Bonferroni Method. The reason for choosing to adjust the p-values is that repeated experiments can increase the chances of accepting the alternative hypothesis so p-values are adjusted according to the number of experiments performed. 10 experiments are performed per sensor.

For the Accelerometer and Proximity, we can see that none of the pairwise groups have a significant difference from each other. This means that even

6. PRE-SURVEY AND EXPERIMENT FINDINGS

Table 6.7: Dunn's Test 1

Groups	Accelerometer	Location	Proximity	Microphone	Camera	Bluetooth
(1,2)	1.0000	1.0000	0.9992	0.8365	1.0000	1.0000
(1,3)	0.4207	0.2084	0.0699	0.0365	0.0732	0.7825
(1,4)	0.0595	0.0125	0.0513	0.0012	0.0048	0.6442
(1,5)	0.2054	0.0010	0.1191	0.0009	0.0007	0.0029
(2,3)	0.6548	0.1774	0.3713	0.8927	0.1694	0.5921
(2,4)	0.1185	0.0077	0.2617	0.2287	0.0123	0.4270
(2,5)	0.3659	0.0005	0.5184	0.1597	0.0018	0.0007
(3,4)	0.8989	0.7642	1.0000	0.8052	0.9040	1.0000
(3,5)	0.9997	0.0869	1.0000	0.6390	0.2947	0.0066
(4,5)	0.9998	0.8360	1.0000	1.0000	0.9617	0.0059

tough the groups perceive sensors differently in general, they all view Accelerometers and Proximity in a similar way.

For the location sensor, it can be observed that groups (1,4), (1,5), (2,4), (2,5) have a significant difference. This can be attributed to the fact that since the groups are formed from the perception of people of the Sensors Feature, the difference in perception between group 1 and group 5 will be larger than between group 1, group 2 and group 1, group 3 since they are not much apart in the scale.

For the microphone sensor, it can be seen that groups (1,3), (1,4) and (1,5) are significantly different from each other. This goes to show that if people rate sensors as even a little intrusive, they all rate the microphone's in a significantly different way than the people who rate sensors as non-intrusive.

For the camera sensor, it can be observed that groups (1,4), (1,5), (2,4), (2,5) have a significant difference in their perception of the intrusion. Similar to the location sensor, people with perception of sensors in general with a lower intrusion level have significantly different responses to the camera intrusion than the people who rate sensors with more intrusion.

For the Bluetooth sensor, there is a significant difference between groups (1,5), (2,5), (3,5) and (4,5). This shows that responses by people who find sensors extremely intrusive is different from the rest of the groups.

Perception of Individual Stakeholders Grouped on the Intrusion of Stakeholders in General

In this section, we try to see if the intrusion level perception by people of Stakeholders in general and the intrusion of the individual stakeholders are related. We try to examine significant differences between groups formed

6.2. Pre-Survey Methodology and Findings

by using question 12's responses on the perception of each stakeholder, and since there are 5 different responses this give five independent groups. The groups 1 to 5 have 7, 11, 32, 69 and 70 people in each respectively. More detailed information about the employment, education, gender and age distribution in the groups is given in tables 6.8, 6.11, 6.9 and 6.10. The mean and variances of each group formed is depicted in figures 6.5a and 6.5b. To start with we examine all the groups simultaneously for all the stakeholder's, we perform the Kruskal-Wallis H test since the data is discrete and not normally distributed. The null hypothesis is that the groups rate the intrusion of a particular stakeholder in a similar way. The alternative hypothesis is that the groups rate the intrusion of a particular stakeholder in a significantly different way. The resulting p-values of this test is displayed in table ???. As it can be seen, the test pronounces that all the groups are significantly different at an alpha with 0.05.

Table 6.8: Employment Classification of Groups

Occupation	1	2	3	4	5
Employed full time	3.92%	4.90%	16.67%	33.33%	41.18%
Employed part time	16.67%	16.67%	8.33%	50.00%	8.33%
Unemployed, looking for work	8.33%	8.33%	16.67%	16.67%	50.00%
Unemployed, not looking for work	0.00%	0.00%	0.00%	33.33%	66.67%
Retired	0.00%	0.00%	0.00%	100.00%	0.00%
Student	4.82%	7.23%	15.66%	37.35%	34.94%
Disabled	0.00%	0.00%	0.00%	0.00%	0.00%

Table 6.9: Gender Classification of Groups

Gender	1	2	3	4	5
Female	5.56%	6.94%	23.61%	36.11%	27.78%
Male	3.97%	6.35%	11.90%	35.71%	42.06%

Table 6.10: Average Birth Year of Groups

1	2	3	4	5
1986	1989	1984	1985	1984

6. PRE-SURVEY AND EXPERIMENT FINDINGS

Table 6.11: Education Classification of Groups

Education	1	2	3	4	5
Less than high school	20.00%	20.00%	0.00%	60.00%	0.00%
High school	15.79%	5.26%	21.05%	31.58%	26.32%
Some college	0.00%	10.00%	20.00%	40.00%	30.00%
Bachelors degree	1.75%	12.28%	15.79%	35.09%	35.09%
Masters degree	5.06%	1.27%	13.92%	37.97%	41.77%
PhD degree	0.00%	7.14%	21.43%	28.57%	42.86%

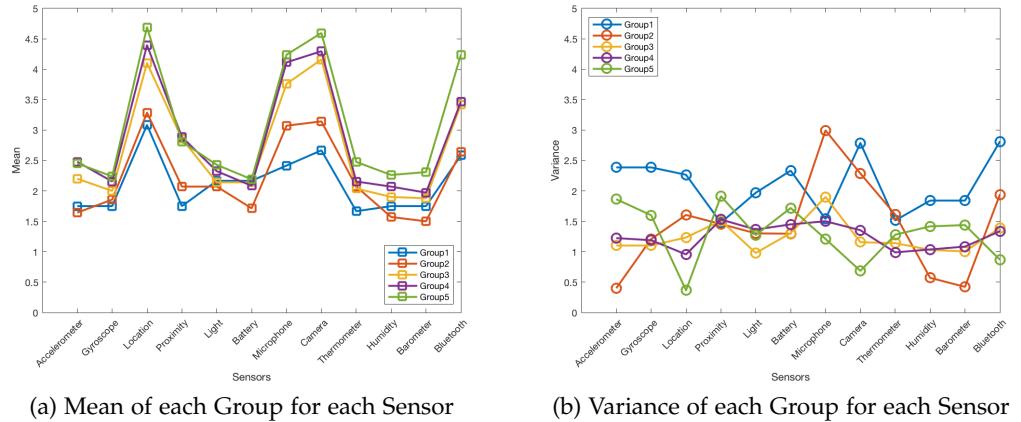


Figure 6.5: Table Schemas

Table 6.12: Kuskal-Wallis Test

Stakeholder	p-value
Corporation	2.1432e-05
Non-Governmental Organization	0.0221
Educational Institution	0.0396
Government	0.0024

This prompts us to take a closer look at which of the groups are significantly different from each other for each stakeholder. For this we continue the experiment with Dunn's Test with p-values adjusted by the Bonferroni Method. The results from the test is shown in table 6.13.

The test was done for all stakeholders since the Kruskal-Wallis test denoted that all groups differ significantly for all stakeholders. Looking at the stakeholder corporation, we see that groups (1,4), (1,5) and (2,5). This goes to show that groups with larger difference in their outlook to stakeholders as

6.2. Pre-Survey Methodology and Findings

Table 6.13: Dunn's Test 1

Groups	Corporation	Non-Governmental Organization	Educational Institution	Government
(1,2)	0.9839	0.8042	0.9565	0.6615
(1,3)	0.3282	0.2351	0.8540	0.0986
(1,4)	0.0240	0.1962	0.6867	0.0289
(1,5)	0.0012	0.0243	0.1254	0.0028
(2,3)	0.9467	0.9992	1.0000	0.9958
(2,4)	0.2047	0.9991	1.0000	0.9252
(2,5)	0.0110	0.7192	0.8415	0.3670
(3,4)	0.6893	1.0000	1.0000	0.9999
(3,5)	0.0197	0.8906	0.4112	0.5794
(4,5)	0.4640	0.6027	0.3427	0.7378

a whole view corporations in a significantly different way.

For the Non-Governmental Organization, only the groups (1,5) differ significantly. This goes to show that groups that do not find stakeholders intrusive and groups that find stakeholders very intrusive rate the intrusion of Non-Governmental Organization in significantly different ways.

For Educational Institutions, none of pairwise comparisons have p-values below 0.05. This goes to show that the intrusion of Educational Institutions by all groups does not differ significantly.

Lastly, for the stakeholder Government, the groups (1,4) and (1,5) differ significantly. This goes to show that groups that view the intrusion of stakeholders with a larger difference view Government in a significantly different way.

The trend observed above is that there is a significant difference in the outlook of individual stakeholders between groups with larger differences in their outlook to stakeholders as a whole, with the exception of Educational Institution where the alternative hypothesis was rejected.

Perception of Individual Contexts Grouped on the Intrusion of Contexts in General

In this section, we examine the relationship between the intrusion of Contexts in General and the individual contexts in question 13. To do this, like the above sections we partition the data into groups based on the answers given questions 14, which asks the user the perception of intrusion of Contexts in general. There are five groups in total. Group one to five have each 12, 11, 42, 74 and 50 people respectively. Additional information about

6. PRE-SURVEY AND EXPERIMENT FINDINGS

the groups on employment, education , gender and year of birth is given in tables 6.14, 6.17, 6.15 and 6.16.

Like in the previous cases, since the data is discrete and not normal, we use the Kruskal-Wallis test to compare the groups perceptions on various contexts. The alpha value is considered to be 0.05. The results are presented in table 6.18. As it can be seen, the test says that there is a significant difference between the groups for all sensors.

Table 6.14: Employment Classification of Groups

Occupation	1	2	3	4	5
Employed full time	4.90%	3.92%	23.53%	36.27%	31.37%
Employed part time	16.67%	8.33%	25.00%	25.00%	25.00%
Unemployed, looking for work	8.33%	16.67%	25.00%	25.00%	25.00%
Unemployed, not looking for work	0.00%	0.00%	0.00%	66.67%	33.33%
Retired	0.00%	0.00%	0.00%	100.00%	0.00%
Student	7.23%	6.02%	20.48%	44.58%	21.69%
Disabled	0.00%	0.00%	0.00%	0.00%	0.00%

Table 6.15: Gender Classification of Groups

Gender	1	2	3	4	5
Male	7.14%	6.35%	23.02%	38.10%	25.40%
Female	5.56%	5.56%	19.44%	38.89%	30.56%

Table 6.16: Average Birth Year of Groups

	1	2	3	4	5
	1986	1986	1986	1985	1983

We not perform Dunn's Test as a post hoc test for all the contexts to observe the exact group pairs that might be significantly different. The results are presented in tables 6.19 and 6.20. For the context Education, the groups (2,5) and (3,5) are significantly different from each other. For the context Entertainment, further investigation shows that groups (1,5), (2,5) and (3,5) are significantly different in each others responses. For the context Environment, groups (2,5) and (3,5) are significantly different from each other. For the context Finance, the groups (1,5), (3,5) and (4,5) are significantly different from each other. In the context health, except for groups (1,5) all the other groups

6.2. Pre-Survey Methodology and Findings

Table 6.17: Education Classification of Groups

Education	1	2	3	4	5
Less than high school	20.00%	0.00%	0.00%	80.00%	0.00%
High school	10.53%	10.53%	31.58%	36.84%	10.53%
Some college 1	0.00%	0.00%	40.00%	30.00%	20.00%
Bachelors degree	7.02%	8.77%	24.56%	35.09%	24.56%
Masters degree	6.33%	3.80%	17.72%	40.51%	31.65%
PhD degree	0.00%	7.14%	17.86%	35.71% 3	9.29%

Table 6.18: Kuskal-Wallis Test

Context	p-value
Education	6.4694e-04
Entertainment	1.0660e-04
Environment	1.3079e-04
Finance	0.0021
Health	0.0011
Shopping	5.4227e-05
Social Network	0.0120
Training	1.2071e-05
Transportation	4.9043e-04

Table 6.19: Dunn's Test Part 1

Groups	Education	Entertainment	Environment	Finance	Health
(1,2)	0.99796	0.99982	0.97055	1.0000	0.63908
(1,3)	1.0000	0.99174	1	0.32963	0.076147
(1,4)	0.94148	0.19262	0.86998	0.22077	0.042547
(1,5)	0.11471	0.0056283	0.21553	0.015364	0.00051115
(2,3)	0.9342	1	0.96665	0.31604	0.9999
(2,4)	0.32744	0.76121	0.083389	0.21508	0.99963
(2,5)	0.0082212	0.082023	0.0048936	0.016365	0.50246
(3,4)	0.83617	0.23034	0.10875	1.0000	1.0000
(3,5)	0.0064191	0.00086418	0.0012905	0.65747	0.32713
(4,5)	0.13825	0.28231	0.60016	0.57519	0.21258

are not significantly different from each other. For the context Shopping, the groups (1,5), (3,4) and (3,5) are significantly different from each other. For the context Social Network, none of the groups are significantly different from each other. In the context Training, the groups (1,5),(3,5) and (4,5) are

6. PRE-SURVEY AND EXPERIMENT FINDINGS

Table 6.20: Dunn's Test Part 2

Groups	Shopping	Social Network	Training	Transportation
(1,2)	1.0000	0.99831	1.0000	1.0000
(1,3)	0.99992	0.94767	1.0000	0.9722
(1,4)	0.18309	0.089135	0.90197	0.23809
(1,5)	0.015842	0.10636	0.010906	0.016376
(2,3)	1.0000	1.0000	1	0.98253
(2,4)	0.39426	0.70552	0.99778	0.30509
(2,5)	0.052334	0.7272	0.068368	0.026144
(3,4)	0.038351	0.21259	0.58123	0.51397
(3,5)	0.00050454	0.29374	2.1697e-05	0.012924
(4,5)	0.69535	1.0000	0.0033189	0.55629

significantly different from each other. Finally for the context Transportation, the groups (1,5), (2,5) and (3,5) are significantly different from each other.

This goes to show that groups that perceive contexts in a more different light tend to have different ways of viewing the individual contexts. We do observe that in some cases (2,5) are significantly different, but (1,5) is not. This can be easily attributed to the low number of responses and the noise in the data.

K- Means to Automatically detect clusters in the Population

We attempt to cluster the population that answered the survey based on the following traits :

- The way people rate data collectors (question 11)
- The way people rate sensors (question 9)
-

For each chosen group of features, we need to discover how many clusters there are in this population. Since we do not have any pre-defined labels given to us, we go through the following procedure:

- Perform the K-Means clustering using various distances
 - Squared Euclidean Distance
 - City Block Distance
 - Cosine Distance
- We plot the Silhouette score to see which of the distances do better
- Visualize the clusters in the highest principal components

6.3 Overview of the Experiment Data

6.4 Findings from the Experiment Data

Chapter 7

Conclusion

Appendix A

Appendix

Bibliography

- [1] Alessandro Acquisti and Jens Grossklags. Privacy and rationality in individual decision making. *IEEE Security & Privacy*, 2(2005):24–30, 2005.
- [2] Rebecca Balebako, Florian Schaub, Idris Adjerid, Alessandro Acquisti, and Lorrie Cranor. The impact of timing on the salience of smartphone app privacy notices. In *Proceedings of the 5th Annual ACM CCS Workshop on Security and Privacy in Smartphones and Mobile Devices*, pages 63–74. ACM, 2015.
- [3] AJ Brush, John Krumm, and James Scott. Exploring end user preferences for location obfuscation, location-based services, and the value of location. In *Proceedings of the 12th ACM international conference on Ubiquitous computing*, pages 95–104. ACM, 2010.
- [4] L Jean Camp. State of economics of information security, the. *ISJLP*, 2:189, 2005.
- [5] Haksoo Choi, Supriyo Chakraborty, Zainul M Charbiwala, and Mani B Srivastava. Sensorsafe: a framework for privacy-preserving management of personal sensory information. In *Workshop on Secure Data Management*, pages 85–100. Springer, 2011.
- [6] Delphine Christin. Privacy in mobile participatory sensing: current trends and future challenges. *Journal of Systems and Software*, 116:57–68, 2016.
- [7] Delphine Christin, Christian Büchner, and Niklas Leibecke. What’s the value of your privacy? exploring factors that influence privacy-sensitive contributions to participatory sensing applications. In *Local Computer Networks Workshops (LCN Workshops), 2013 IEEE 38th Conference on*, pages 918–923. IEEE, 2013.

BIBLIOGRAPHY

- [8] Dan Cvrcek, Marek Kumpost, Vashek Matyas, and George Danezis. A study on the value of location privacy. In *Proceedings of the 5th ACM workshop on Privacy in electronic society*, pages 109–118. ACM, 2006.
- [9] George Danezis, Stephen Lewis, and Ross J Anderson. How much is location privacy worth? In *WEIS*, volume 5. Citeseer, 2005.
- [10] Linda Deng and Landon P Cox. Livecompare: grocery bargain hunting through participatory sensing. In *Proceedings of the 10th workshop on Mobile Computing Systems and Applications*, page 4. ACM, 2009.
- [11] Fosca Giannotti, Dino Pedreschi, Alex Pentland, Paul Lukowicz, Donald Kossmann, James Crowley, and Dirk Helbing. A planetary nervous system for social mining and collective awareness. *The European Physical Journal Special Topics*, 214(1):49–75, 2012.
- [12] Eiji Hayashi, Oriana Riva, Karin Strauss, AJ Brush, and Stuart Schechter. Goldilocks and the two mobile devices: going beyond all-or-nothing access to a device’s applications. In *Proceedings of the Eighth Symposium on Usable Privacy and Security*, page 2. ACM, 2012.
- [13] Jialiu Lin, Shahriyar Amini, Jason I Hong, Norman Sadeh, Janne Lindqvist, and Joy Zhang. Expectation and purpose: understanding users’ mental models of mobile app privacy through crowdsourcing. In *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*, pages 501–510. ACM, 2012.
- [14] Emiliano Miluzzo, Nicholas D Lane, Shane B Eisenman, and Andrew T Campbell. Cenceme—injecting sensing presence into social networking applications. In *European Conference on Smart Sensing and Context*, pages 1–28. Springer, 2007.
- [15] Prashanth Mohan, Venkata N Padmanabhan, and Ramachandran Ramjee. Nericell: rich monitoring of road and traffic conditions using mobile smartphones. In *Proceedings of the 6th ACM conference on Embedded network sensor systems*, pages 323–336. ACM, 2008.
- [16] Evangelos Pournaras. Application form to the research ethics committee of eth zurich.
- [17] Evangelos Pournaras, Jovan Nikolic, Pablo Velásquez, Marcello Trovati, Nik Bessis, and Dirk Helbing. Self-regulatory information sharing in participatory social sensing. *EPJ Data Science*, 5(1):1, 2016.
- [18] Ashwini Rao, Florian Schaub, and Norman Sadeh. What do they know about me? contents and concerns of online behavioral profiles. *arXiv preprint arXiv:1506.01675*, 2015.

Bibliography

- [19] Lei Song, Yongcai Wang, Ji-Jiang Yang, and Jianqiang Li. Health sensing by wearable sensors and mobile phones: a survey. In *e-Health Networking, Applications and Services (Healthcom), 2014 IEEE 16th International Conference on*, pages 453–459. IEEE, 2014.
- [20] Jinyan Zang, Krysta Dummit, James Graves, Paul Lisker, and Latanya Sweeney. Who knows what about me? a survey of behind the scenes personal data sharing to third parties by mobile apps. *Proceeding of Technology Science*, 2015.

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

First name(s):

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.