



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	7
3.1 Introduction	7
3.2 Model Intricacies	7
3.2.1 Categorization of the Features	8
3.2.2 Categorization of the Sub-Features	9
3.2.3 Weight Matrix Calculation	10
3.2.4 Cost Matrix Calculation	11
3.2.5 Cost and Privacy Metrics	11
3.2.6 Improving the Metrics	13
3.2.7 Summarization of Collected Data	13
3.3 Analysis of the Model	14
3.3.1 Setup	14
3.3.2 Results	15
4 Experiment Methodology	21
4.1 Preparatory Phase	21
4.1.1 Pre-Survey	21
4.1.2 Sub-Features	22
4.1.3 Privacy Options	22
4.1.4 Question Structure	22
4.1.5 Budget and Experiment Duration	23
4.2 Entry Phase	23
4.2.1 Collecting General User Information	23
4.2.2 Categorization of Features	23

CONTENTS

4.2.3	Categorization of Sub-Features	24
4.2.4	Answering Questions with No Incentives	26
4.3	Core Phase	28
4.3.1	Improve Privacy or Credit	30
4.3.2	Answering Questions with Incentives	30
4.4	Exit Phase	32
4.5	FairDataShare Web Portal	32
4.5.1	Data Generator's Portal	32
4.5.2	Stakeholder's Portal	34
5	Explanation of the Mobile Application	37
5.1	The Building Blocks	37
5.2	The Mobile Application	37
5.2.1	Local Storage	37
5.2.2	Alarms	37
5.2.3	Privacy and Credit Improvement	37
5.2.4	Recommendations	37
5.2.5	Recording User Choices	37
5.2.6	Sensor Data Collection and Summarization	37
5.2.7	Server Synchronization	37
5.3	The Server	37
5.3.1	Kinvey Data Storage	37
5.3.2	FairDataShare Web Portal	39
6	Pre-Survey and Experiment Findings	41
6.1	Overview of the Pre-Survey Data	41
6.2	Pre-Survey Methodology and Findings	41
6.3	Overview of the Experiment Data	41
6.4	Findings from the Experiment Data	41
7	Conclusion	43
A	Appendix	45

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

Chapter 3

Computational Model

3.1 Introduction

Why the model is needed, add some previous paper about incentives and what we do differently Our aim is to create a computational model that is able to collect useful data about the influence of incentives on mobile data sharing. (Quote some studies that have done similar studies with no data incentives). In the model, what we try to do is first create a user profile by asking the user some preliminary questions. Then, we proceed to assign each sensor data request with a maximum achievable credit using the formulated user profile. The idea of the model is that it assigns requests where the user least desires to share mobile sensor data with higher incentive costs. Respectively, the data requests where the user desires to share data more are assigned lower incentive costs. This permits us to see whether incentives do indeed make a difference in data sharing, since assigning more credit to data requests where the user would anyway give mobile sensor data is futile to our goals. The model aims to collect useful data that examine the relationship between incentives and mobile data sharing, mainly when the user least desires to share data.

3.2 Model Intricacies

The sections below explain the various building blocks of the computational model. The Figure 3.1 provides an overview of how the model works. To begin with the model, each user is asked to enter various non-intrusive personal information that can help in analysing the user's behavior. For example, this can consist of the age, gender, country of residence and employment.

3. COMPUTATIONAL MODEL

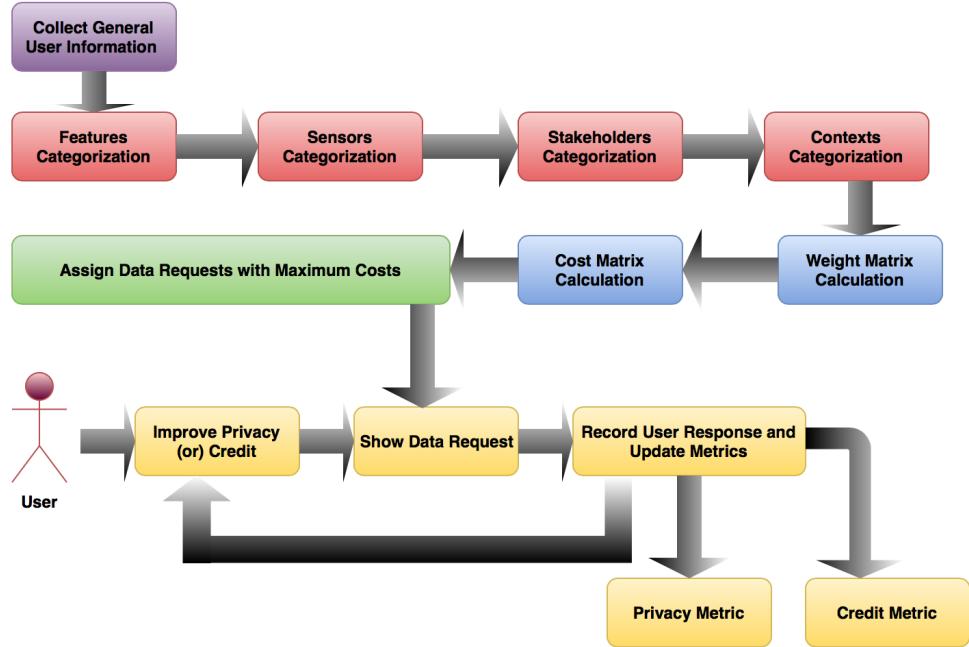


Figure 3.1: Computational Model Flow Chart

3.2.1 Categorization of the Features

After we are done recording the user's personal information, we go into the categorization of the features. In this model, features consist of Sensors, Stakeholders and Contexts. The Sensors consist of the sensors in the mobile phone which the user's can trade. Let the category assigned to the Sensors be represented by S . The Stakeholders consist of any entity that can request the user for mobile sensor data. Let the category assigned to the Sensors be represented by DC . The Contexts consist of the purpose for which a Stakeholder would like to obtain the user's mobile sensor data. Let the category assigned to the Sensors be represented by C . Categorization of the features means that the user places each of the features into predefined categories, and 2 or more features can be placed into the same category. The user is asked to categorize each feature in the available cat number of categories. The first category indicates that the feature does not contribute much to the data sharing decision. Respectively, the category cat indicates that the feature contributes a lot to the data sharing decision. The categories are linearly scaled and equally spaced. The reason categorization was chosen is as to not rule out the possibility to that two features may be considered equally important in the data sharing decision and this may be missed by ranking the features. Once the user has categorized the Sensors, Stakeholders and the Contexts, the weights of each feature in the data sharing decision is calcu-

lated. The category feature Sensors has been placed into be represented by S , the category feature Stakeholders has been placed into be represented by DC and the category feature Context has been placed into be represented by C . Hence the respective weights $weights_S$, $weight_{DC}$ and $weight_C$ are calculated as follows :

$$weights = \frac{S}{S + DC + C} \quad (3.1)$$

$$weight_{DC} = \frac{DC}{S + DC + C} \quad (3.2)$$

$$weight_C = \frac{C}{S + DC + C} \quad (3.3)$$

3.2.2 Categorization of the Sub-Features

Once the features have been categorized and their weights calculated, the sub-features need to be categorized. In this model, sub-features consist of the various types of Sensors available on the mobile phone, the various types Stakeholders that request mobile data from users and the different types of Contexts for which mobile data is requested. In other words, sub-features are the different kinds of features that appear during data request to the user. The following are examples of sub-features for each feature :

- Sensors : Accelerometer, Battery and Gyroscope
- Stakeholders : Company, Non-Governmental Organization and Government.
- Contexts : Education, Entertainment and Navigation.

For each of the features, the respective sub-features are assigned a unique identifier ranging from one to the length of sub features. Now for each of the features, the respective sub-features need to be in turn categorized in a similar fashion to section 3.2.1. Each sub-feature can be placed in the available cat categories. The first category indicates that the user finds the sub-feature very non privacy intrusive. This means that the user would not be worried trading data for this sub-feature. The last category indicates that the user finds this sub-feature very privacy intrusive. This in turn means that the user would be reluctant giving data for this sub-feature. All the categories in between are linearly scaled and equally spaced. The user then places for each feature, the respective sub-features in the cat available categories according to the perceived intrusion level. A conceptual diagram is shown in figure 3.2. For the sub-features of Sensors, categories they are placed in are represented by S_i , where S is the feature Sensors and i is the id of the

3. COMPUTATIONAL MODEL

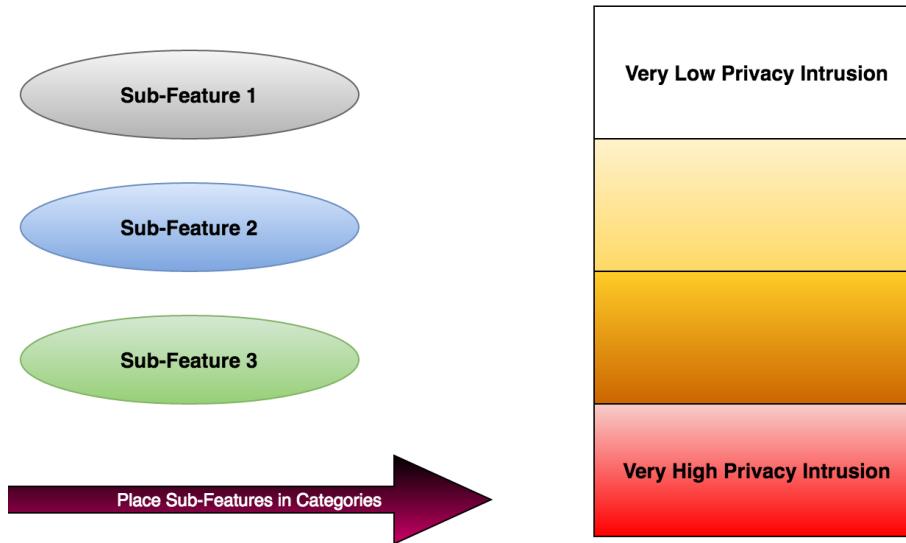


Figure 3.2: Categorizing Sub-Features according to the perceived Intrusion Level

sub-feature. Similarly, categories assigned to sub features of features Stakeholders and Contexts respectively are represented by DC_j and C_k , where j and k are the id's of the sub-features categorized.

3.2.3 Weight Matrix Calculation

Each data request to the user consists of the 3 features in them. Each of those features has a number of sub-features that can appear in turns in a data request, that is in a factorial form. Let $count(feature)$ be function that gives the number of sub-features given a feature. The total number of data requests is :

$$N_{DR} = count(Sensors) * count(Stakeholders) * count(Contexts) \quad (3.4)$$

Let WM be a matrix with dimensions $count(Sensors) \times count(Stakeholders) \times count(Contexts)$. We call this the weight matrix. The cell $WM_{i,j,k}$ represents 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 . That is, each cell of WM represents a data request to the user. The aim of the weight matrix is to use the information collected from the user profiling to assign various weights to each data requests. Intuitively, the process examines the data requests where the user is least likely to trade data and assigns higher weights to those data requests. This process can be seen in section 3.3 in more detail. 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 :

$$WM_{i,j,k} = (S * S_i) + (DC * DC_j) + (C * C_k) \quad (3.5)$$

Applying this formula to every cell gives the weight matrix WM .

3.2.4 Cost Matrix Calculation

Now that the weights for every data request has been calculated, we need to calculate the exact amount of money the users can receive for a data request. Let CM be the cost matrix with dimensions $count(Sensors) \times count(Stakeholders) \times count(Contexts)$. Let us assume to have a budget of B for a day, where B can be in an actual currency or any sorts of virtual credits. For now, the budget will be referred to as 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 :

$$sum(WM) = \sum_{i=1}^{count(Sensors)} \sum_{j=1}^{count(Stakeholders)} \sum_{k=1}^{count(Contexts)} w_{i,j,k} \quad (3.6)$$

where the function $sum(matrix)$ 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 :

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

Doing this for every cell of CM , the whole cost matrix can be calculated. Now, we have the credits allocated per day for every data request.

3.2.5 Cost and Privacy Metrics

Every data request now has an associated 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 for one day. Similarly, the Privacy metric is the amount of privacy percentage the user has maintained. That is, it intuitively quantifies the amount of data the user has refused to share hence implying privacy. The Cost and Privacy are inversely proportional to each other, in the sense that when the Cost goes up and Privacy goes down

3. COMPUTATIONAL MODEL

and vice versa. For each data request, the user can choose how much data is to be shared, from the maximum amount of data to no data at all. Each option corresponds to a summarization level explained in detail in section 3.2.7. 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 where the user gives all the data requested and m to where the user chooses not give any data at all. Therefore there are a total of m options for a data request. While assigning costs there are two scenarios:

- Assigning option costs without a participation cost.
- Assigning option costs inclusive of a participation cost.

Let us examine the first scenario. Let us assume that we are calculating the option costs for data request with Sensors sub-feature i , stakeholders sub-feature j and contexts sub-feature k . Let us calculate the assigned cost for option number h of this data request:

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

Applying this formula by replacing h by the options from 1 to m gives the cost the user receives for each option. Similarly, if you would like to assign a participation cost to each option, it would mean that even though the user does not share data, they still receive some money for answering the data request. This concept can be implemented to ensure user participation. (Quote some paper 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 cost $frac_h$ an option numbered h gets:

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

Now that we know the fraction of the cost option f can be assigned, 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 arithmetic average of all the costs and privacy obtained from every answered data request. If a data request is left unanswered, maximum privacy and minimum cost is assumed.

3.2.6 Improving the Metrics

Before the user 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.

3.2.7 Summarization of Collected Data

As mentioned before, each data request can have options m number of options the user can choose from. 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 than in its original form. The lower the summarization level gives data closer to its original form. In this model, 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 whole day, and at the end of the day according to the option chosen by the user, it is summarized. Summarization can be linearly assigned to each option starting with the highest privacy corresponding to highest summarization level , that is no data sharing to the lowest summarization level, that is no summarization at all. An example of assigning the summarization level $summ_h$ for option h 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. COMPUTATIONAL MODEL

3.3 Analysis of the Model

In this section, we take a scenario of the computational model and show how exactly the model works. In particular, the focus is on how the model varies the weights to questions according to the user input.

3.3.1 Setup

the sensors, stakeholders, and contexts and other special parameters such as number of options and all To explain the model using examples, we take into consideration the following sub-features for each feature:

1. Sensors
 - a) Accelerometer -1
 - b) Noise -2
 - c) Location -3
2. Stakeholders
 - a) Corporation -1
 - b) Government -2
 - c) Educational Institution -3
3. Contexts
 - a) Navigation -1
 - b) Environment -2
 - c) Social Media -3

The numbers indicated next to the sub-features is the sub-feature identifier. This uniquely identifies a sub-feature within a feature category. Each user will receive an amount of

$$\text{count}(\text{Sensors}) * \text{count}(\text{Stakeholders}) * \text{count}(\text{Contexts}) = 27$$

data requests in total. Each data request has five privacy options ranging from one to five. the option one indicates the users would like to trade all their data, and option five indicates the users refuse to share data their for this data request. Additionally, it is assumed that the core phase has a Budget $B = 100$ per day. The input to the model are the user choices during the categorization of the features and sub-features.

3.3.2 Results

In this section, three user scenarios will be introduced and explained in order to explore the properties of the weight and cost matrices. First, we will begin by introducing the way the user has categorized the features and sub-features. 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. For example, accelerometer is Sensor's sub-feature 1. Similarly, Navigation is Context's sub-feature 1. 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

In scenario 1.1, 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 receive 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, the formulation of the weight matrix can be seen in figure 3.3a, and the cost matrix can be seen in figure 3.3b. As we expected, for each data request indicated as a tuple of (sensors, stakeholders, contexts) in the x-axis of figures 3.3 have identical weights and costs. This is due to the fact that the user finds all the Features and sub-features equally intrusive so all the data requests are weighted equally.

Table 3.1: Categorization for Scenario 1.1

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

The theory that all equally intrusive Features and sub-features should have data requests with equal weights and costs forms the basis of the computa-

3. COMPUTATIONAL MODEL

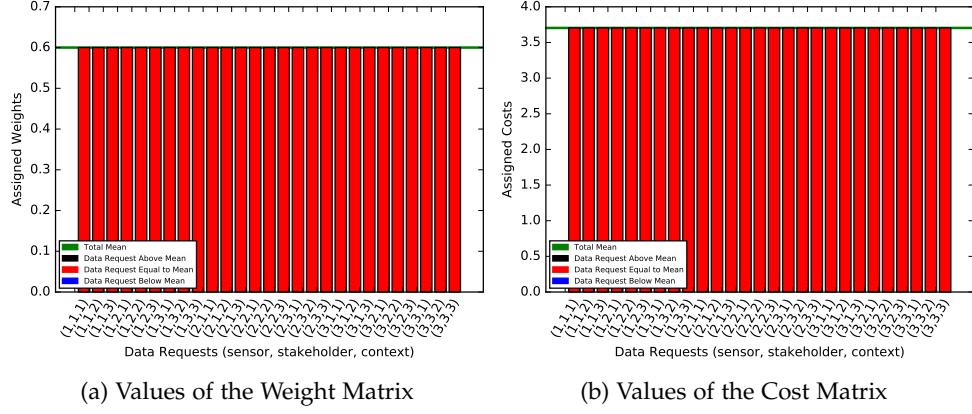


Figure 3.3: Examining Scenario 1.1

tional model. Hence, it becomes essential to view another similar scenario to confirm that this indeed works. The table 3.2 is the user input to the next scenario 1.2. Similar to scenario 1.1 but with different inputs, the Features and sub-features categorized are viewed to all be equally intrusive by the user. Hence as shown in figures 3.4a and 3.4b, data requests are again weighted equally.

Table 3.2: Categorization for Scenario 1.2

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

We can conclude that if the user perceives the feature and respective sub-features in an equally intrusive way, then all the data requests will have the same costs assigned.

Scenario 3

Table 3.3 indicates the user input to scenario 3. 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 Sensor's sub-features with identifiers 1,2 and 3 have respectively categories 1,3 and 5. This means that

3.3. Analysis of the Model

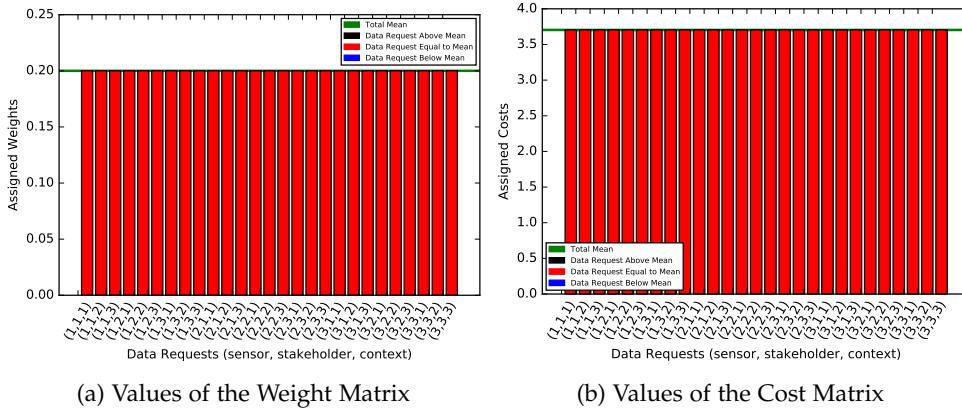


Figure 3.4: Examining Scenario 1.2

requests requests with Sensor's sub-feature 1 will have the lesser weight in comparison to the other Sensor's sub-features. Similarly, the data requests with Sensor's sub-feature 2 will have a higher weightage than Sensor's sub-feature 1, but lesser than Sensor's sub-feature 3. Lastly, data requests with Sensor's sub-feature 3 will have a higher weight compared to the others, due to its category being 5. The weight and cost matrices can be seen in figures 3.5a and 3.5b 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
3	1	3	5
Stakeholders	Corporation	Government	Educational Institution
3	3	3	3
Contexts	Navigation	Environment	Social Media
3	3	3	3

From the above input and graphs, we can conclude that the model assigns a higher weight to data requests with sub-features that the user finds more intrusive compared to the others.

Scenario 4

An attempt is made to vary the feature and sub-feature categories at once, to show how varying their values together affects the assignments of the weight matrix. Table 3.4 is the user input to the scenario 4. All the Features

3. COMPUTATIONAL MODEL

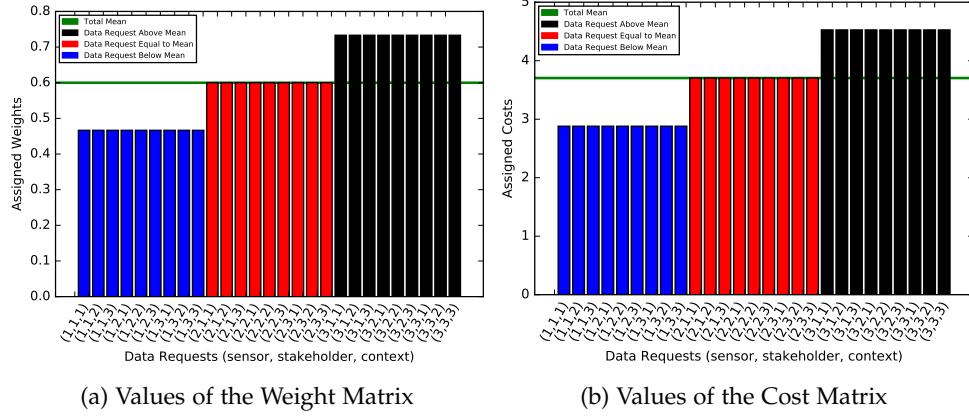


Figure 3.5: Examining Scenario 3

have different categories assigned from 3 to 5. Additionally, the sub-feature 1 of each feature has a category of 5, higher than the others which are all categorized as 1. The weight and cost matrices generated for this scenario can be seen in figures 3.6a and 3.6b respectively.

As it is observed for 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 all sub-features 1 of each feature. It happens because all of the sub-features 1 are assigned a category of 5. The feature Sensors feature 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 3.6, 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. 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 request 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 Context feature's sub-feature 1 has a category of 5. This is due to the fact that Sensor's feature and Stakeholders feature have a higher categories of 5 and 4 respectively than the context 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 its feature.

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

3.3. Analysis of the Model

themselves. As seen in the figure, all the features with higher intrusion categorizations have weights and costs that are well above average.

We can conclude 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 with a sub-feature lowly categorized. Features with lower categories contribute lesser to the weight assignments, irrespective of their sub-feature categories.

Table 3.4: Categorization for Scenario 4

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

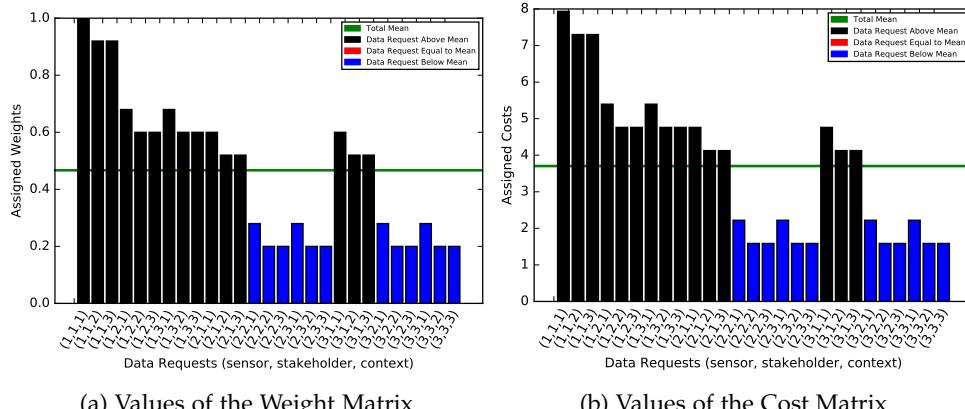


Figure 3.6: Examining Scenario 4

Chapter 4

Experiment Methodology

4.1 Preparatory Phase

4.1.1 Pre-Survey

(For each sensor, dc and context show graphs from the pre-survey why each of them was chosen)

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:

1. Sensors - The mobile mobile sensor data shared
2. Stakeholders - The entity to whom mobile sensor data is shared
3. Context - The purpose for which mobile sensor data is shared

From each of the above mentioned features, there were a lot of sub-features to choose from. Sub-features are the elements that come under Feature. For example, Light sensor is a sub-feature of Sensors and Corporation is a sub-feature of Stakeholders. 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 the above(link to appendix). Also, it can help us redesign some of the aspect 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 privacy and sensors. Participants were not paid for filling out the survey.

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

4. EXPERIMENT METHODOLOGY

4.1.2 Sub-Features

194 pre-survey entries have been filled out by the participants. Using the data obtained from the pre-survey, four sub-features were chosen from each feature.

4.1.3 Privacy Options

Each data request is accompanied with privacy options ranging from 1 to 5. Option 1 indicates that the users would like to share their raw data without any sort of summarization or filter. 5 indicates that the users would not like to share 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 refer to 3.2.7.

4.1.4 Question Structure

A data request is when a stakeholder asks users for mobile sensor data for a particular context or purpose. From now on, we will refer to mobile sensor data as just data. 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
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.

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, 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

explanation of screen shots

4.2.1 Collecting General User Information

As the figure 4.1 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 submit button is pressed in the screen 4.1c, 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 need to categorize the features Sensors, Stakeholders and Contexts. As shown in figure 4.2a, the features are indicated followed by 5 options of privacy ranging from "very low privacy intrusion" to "very privacy high privacy intrusion". The option "very low

4. EXPERIMENT METHODOLOGY

The figure consists of three screenshots of a mobile application interface titled "GetUserInformation".

- Screen 1:** Employment Status. A grid of radio buttons for employment status: Full Time, Part Time, Not Looking for Work, Looking for Work, Retired, Student, and Disabled. The "Student" button is selected.
- Screen 2:** In what country did you spend most of your life? A dropdown menu with "France" selected. A green "SUBMIT" button is visible below it.
- Screen 3:** How often do you check your mobile phone a day? A grid of radio buttons for frequency: <35, 36-70, 71-100, 101-130, and >130. The "71-100" button is selected.

(a) User Information Screen 1 (b) User Information Screen 2 (c) User Information Screen 3

Figure 4.1: User Information Screens

"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 options are compulsory, and no default option is provided. Users cannot navigate to the next page without filling out all the questions.

4.2.3 Categorization of Sub-Features

For each of the features categorized in the previous sub-section, each of the sub-features need to be categorized in a similar fashion. Again, the privacy options range from very low privacy intrusion" to "very privacy high privacy intrusion" like in section 4.2.2 . The users are first presented with the

4.2. Entry Phase

The figure consists of two side-by-side screenshots of a mobile application. Both screenshots show a top navigation bar with various icons and the time '23:29' or '23:30'. Below the navigation bar, there are two main sections: 'Categorize Features' on the left and 'Categorize Sensors' on the right.

(a) Categorizing Features:

- Section title: 'Categorize Features'
- Text: 'How intrusive are the following features of information sharing:'
- Category: 'Sensors'
- Rating: 'very high privacy intrusion' (highlighted in orange)
- Category: 'Data Collectors'
- Rating: 'medium privacy intrusion' (highlighted in orange)
- Category: 'Context / Purpose'
- Rating: 'very low privacy intrusion' (highlighted in orange)
- Bottom button: 'SUBMIT' (highlighted in green)

(b) Categorizing Sensors:

- Section title: 'Categorize Sensors'
- Text: 'How intrusive are the following sensors of information sharing:'
- Category: 'Accelerometer'
- Rating: 'medium privacy intrusion' (highlighted in orange)
- Category: 'Location'
- Rating: 'very high privacy intrusion' (highlighted in orange)
- Category: 'Light'
- Rating: 'low privacy intrusion' (highlighted in orange)
- Category: 'Noise'
- Rating: 'high privacy intrusion' (highlighted in orange)
- Bottom button: 'SUBMIT' (highlighted in green)

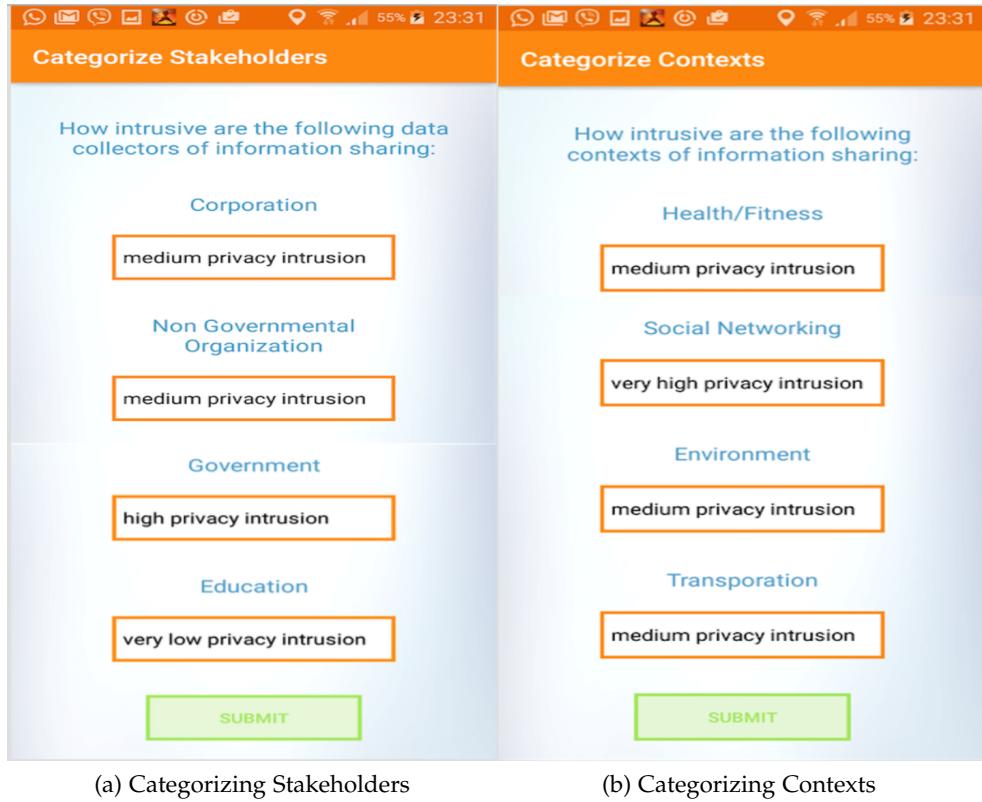
Figure 4.2: Categorizations

categorization of Sensors sub-features as shown in figure 4.2b. 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 an 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 4.3a.

Each stakeholder type has a drop down menu each where the user can again classify how much each of them affect data sharing. Once the user has finished entering the privacy intrusion level for stakeholders, the user can click the green submit button and is directed to the next page. On this page, the user will need to categorize how much each of the Context's sub-features affect mobile sensor data sharing. This is depicted in figure 4.3b. Each context has a drop down menu below, where the user can rate each context from "low privacy intrusion" to "very privacy high privacy intrusion". 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

4. EXPERIMENT METHODOLOGY



The figure consists of two side-by-side screenshots of a mobile application. Both screenshots show a header with various icons and the time '23:31'.

(a) Categorizing Stakeholders:

Question: How intrusive are the following data collectors of information sharing:

- Corporation: medium privacy intrusion
- Non Governmental Organization: medium privacy intrusion
- Government: high privacy intrusion
- Education: very low privacy intrusion

(b) Categorizing Contexts:

Question: How intrusive are the following contexts of information sharing:

- Health/Fitness: medium privacy intrusion
- Social Networking: very high privacy intrusion
- Environment: medium privacy intrusion
- Transporation: medium privacy intrusion

Both screens have a green 'SUBMIT' button at the bottom.

Figure 4.3: Categorizations

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 recorded, users will be presented with 64 questions. Each question 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 is the sensor data given away for that request and vice versa. Users can change the answers to 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 4.4a. After the users choose an option for the data request, a green submit button appears which is shown in figure 4.4b. After clicking on the submit button, response to the data request cannot be changed. At this time, no indications of credit gain or privacy improvements are indicated.

4.2. Entry Phase

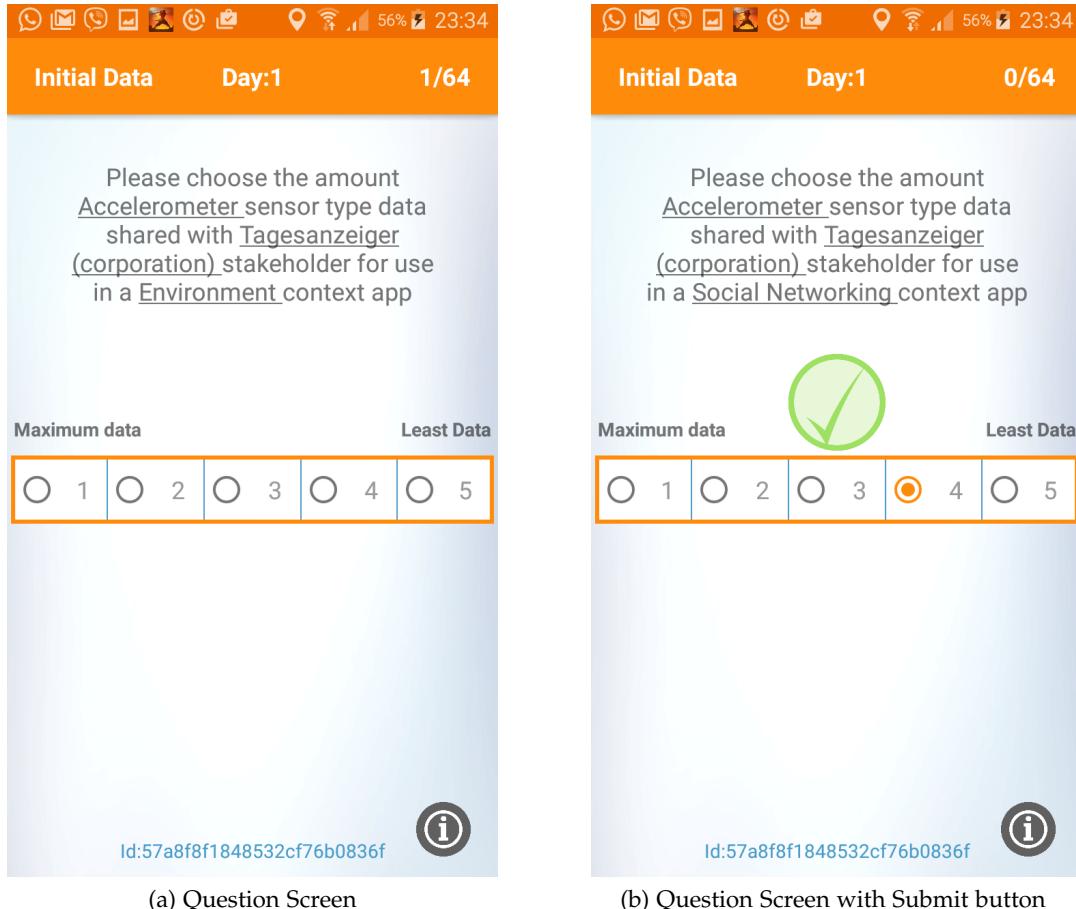


Figure 4.4: First Day Screen

The "*i*" button at the bottom right of the screen is clickable. This takes the user to the FairDataShare portal. Figure 4.5b shows the homepage of the portal. Users need to then click on the data generator registration section of the website where users can signup with their:

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

The unique identifier is located at the bottom of the page is an alphanumeric sequence. If it is long pressed the user can select the identifier, then copy and paste it in the textbox asking for the unique identifier. Figure 4.5b shows what the registration page looks like. The users can use this website to see

4. EXPERIMENT METHODOLOGY

all the data collected from them for all the mobile sensors. More details about the FairDataShare portal refer to the section 4.5.

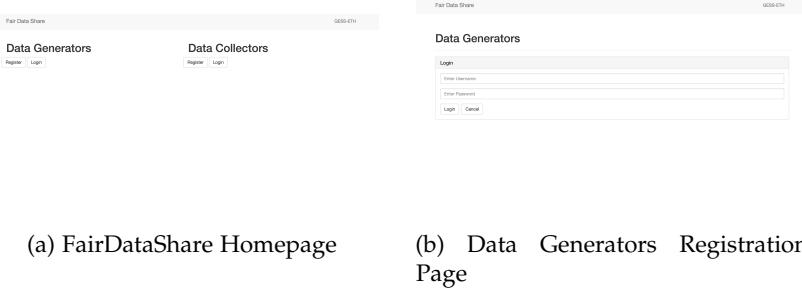


Figure 4.5: FairDataShare Portal

In the task-bar, the user can see the bidding day number and how many questions have been answered from the total available. Day number one corresponds to the day where users answer questions with no incentives of any kind. Once all the questions have been answered, the user goes to the core phase of the experiment, which starts at day number two.

4.3 Core Phase

Once the entry phase is done, the user is presented with the screen shown in figure ???. The first presented screen 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 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. The button numbered 9 is the button that takes the user to the FairDataShare portal. 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 stakeholders. The item numbered 8 is the unique identifier of the user. This can be selected and copied by long pressing the unique identifier for one second. 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 when all the 64 have been answered, the number of questions answered is reset and the number of round 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.

4.3.1 Improve Privacy or Credit

The improvement screen shown in figure 4.6 is where users can choose whether he would like to improve the privacy or the credit that has been obtained. 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 credit already obtained further. This uses algorithm that uses the previous user answers to put forth a data request that can increase the credit to the maximum. The credit improvement button is represented by the 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. Then again, 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 example for each button is given in the next section after introducing the next screen.

For example, if a user chooses to improve the privacy, then clicks on improve privacy button and gets a data request, but still chooses option maximum data (least privacy) for that 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. Trading all data gives the user more credit, but decreases the privacy metric.

Similarly, if a user chooses to improve the credit, then clicks on the improve credit button and gets a data request. Then the user chooses the option least data (maximum privacy) which indicates that no data is traded for this request, this is counters the initial desire to improve the credit obtained.

4. EXPERIMENT METHODOLOGY

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 4.7a. This screen is called the "bidding screen". This screen is very similar to the screen 4.6 presented in the entry phase, except that the user is aware of the amount of privacy and credit obtained. Additionally, the user can see information about how the privacy and credit will increase or decrease for each data request, according to the chosen option. The items numbered 11 are the possible answers ranging from one to five (option numbers are not indicated on the screen). The items numbered 12 are the improvement in privacy for each possible option of the current data request. The items numbered 13 are the improvements in credit for each possible options of the current question. 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 ?? and then click again on the green submit button shown in 4.7b 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 in figure 4.8a. 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 3 to a data request and has clicked on improve credit, the system puts an orange box around options 1,2 and 3. Similarly, if the user clicked on improve privacy button, the system would recommend the options 3,4 and 5. Two examples of this are provided in figures 4.8.

4.4 Exit Phase

After the end of the core phase, the participants are asked to fill up a survey based on their experience of 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 survey ² 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.

4.5 FairDataShare Web Portal

The FairDataShare portal ³ 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

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

Once the users are registered, they can come back to the portal after 24 hours period or later to view their mobile sensor data collected in the server. The data portal login page is shown in figure 4.9a. 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 done, users will be redirected to the data collection page shown in figure 4.9b with the following options in the task-bar :

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

²https://descil.eu.qualtrics.com/SE/?SID=SV_3P0ySMqNe006v5j

³<http://fair-data-share.inn.ac/>

4. EXPERIMENT METHODOLOGY

Figures 4.10a, 4.10b, 4.11a and 4.11b show examples of data that can be seen for the location, light, accelerometer and noise sensor.

In the experiment, day number one is the entry phase, the core phase is day number two and three.

4.5.2 Stakeholder's Portal

For a stakeholder to view data, they need to register in the portal shown in figure 4.5a by clicking register. Once that is done, the page in figure ?? is shown asking for :

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

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

1. Corporation
2. Educational Institution
3. Government
4. Non-Governmental Organization

After this, the stakeholder can click on the register button. Once registered, the stakeholder can login like in 4.12b. Access is then granted to the page in figure 4.13. The stakeholder can choose from each drop down list:

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

Once this is entered, the stakeholder can see 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 request. Stakeholders can view sensor data in a similar fashion to users shown in figures 4.10 and 4.11.

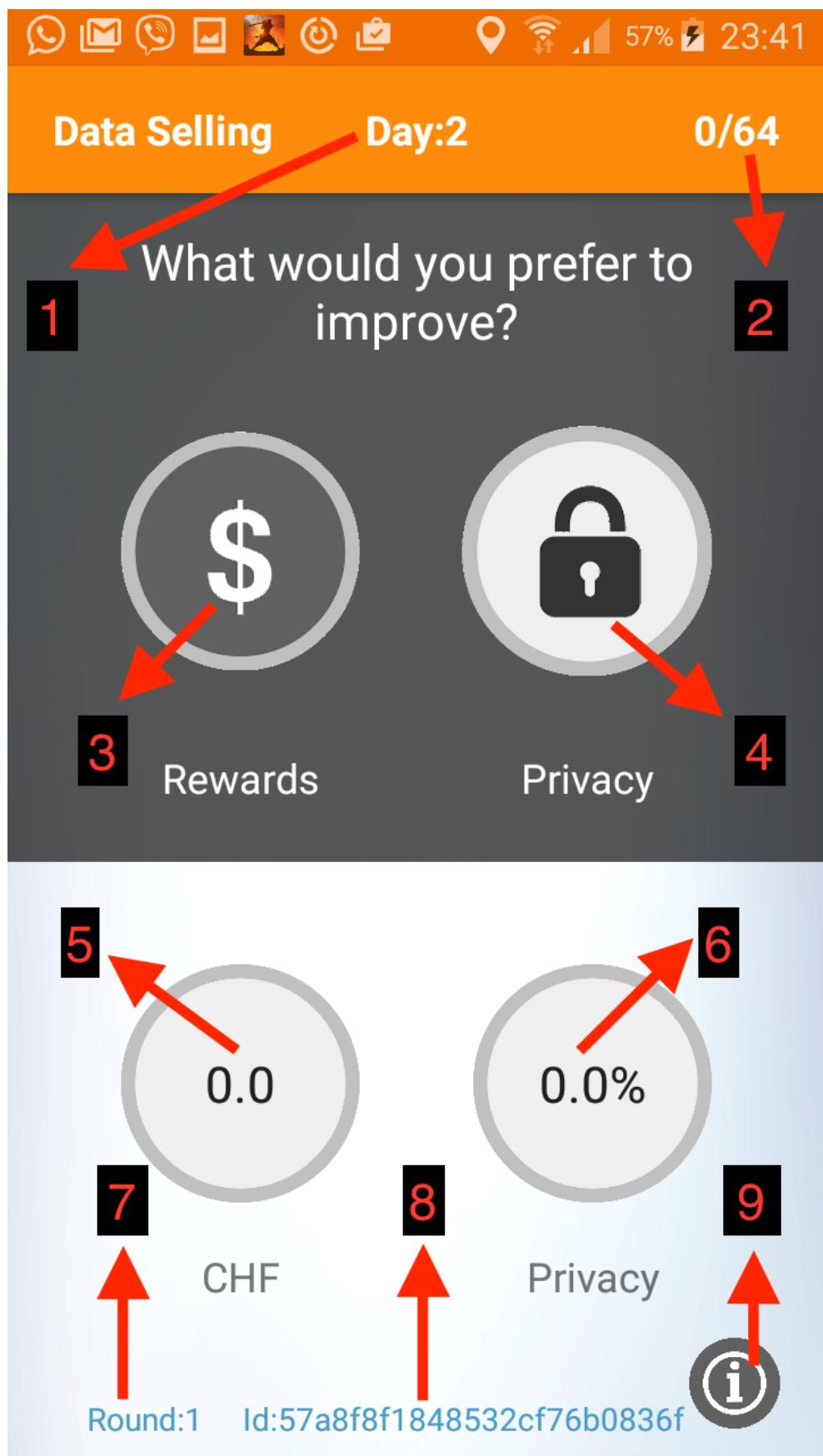


Figure 4.6: Improvement screen

4. EXPERIMENT METHODOLOGY

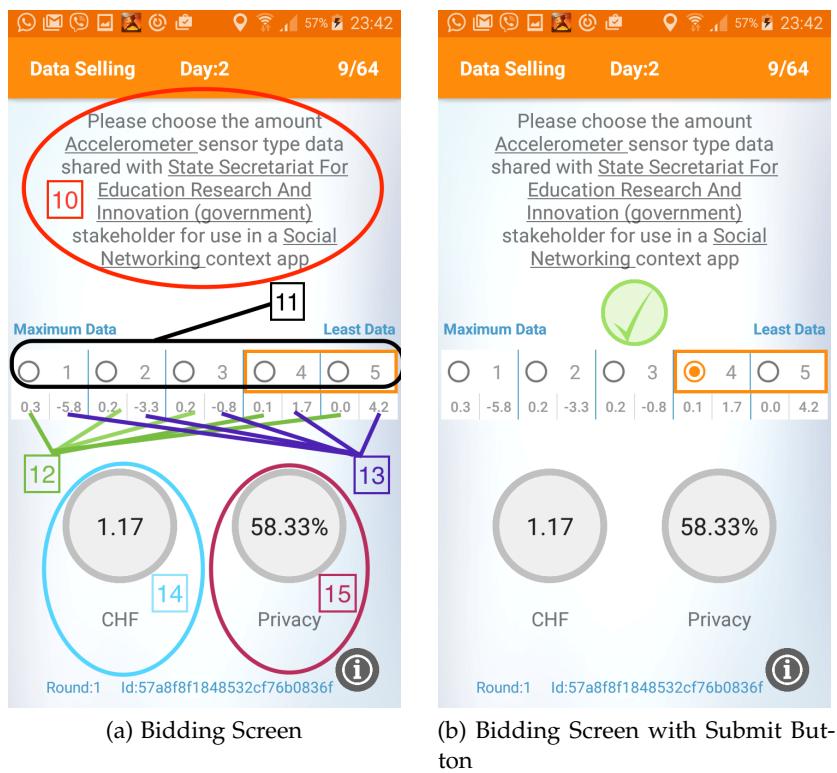


Figure 4.7: FairDataShare Portal

4.5. FairDataShare Web Portal

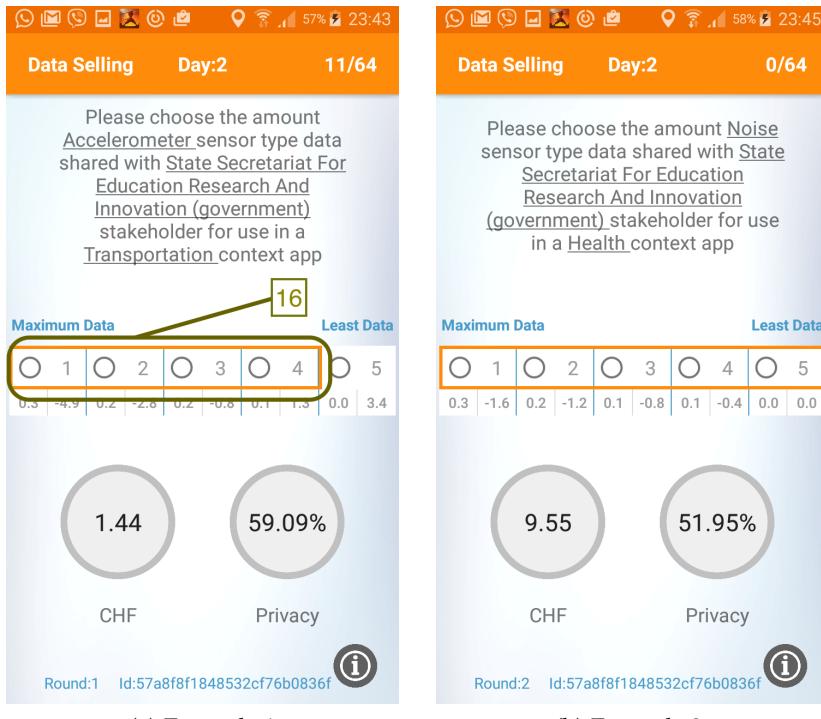


Figure 4.8: Recommendation Box

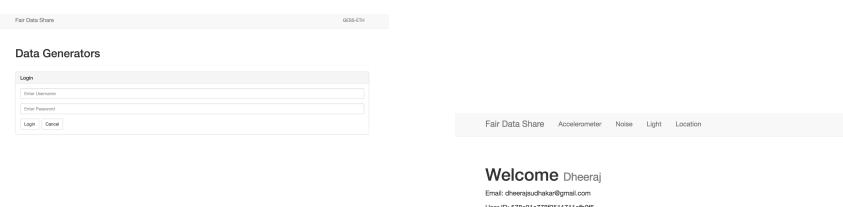


Figure 4.9: Entering the Portal

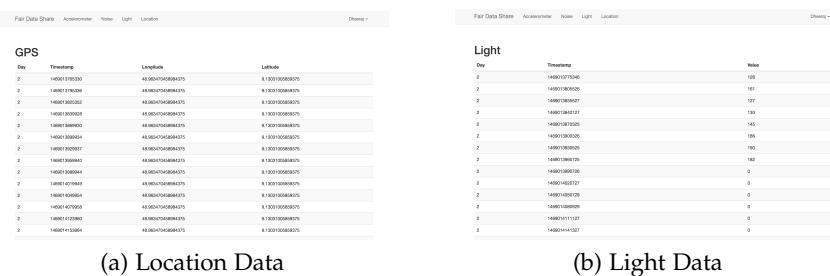


Figure 4.10: User Data

4. EXPERIMENT METHODOLOGY

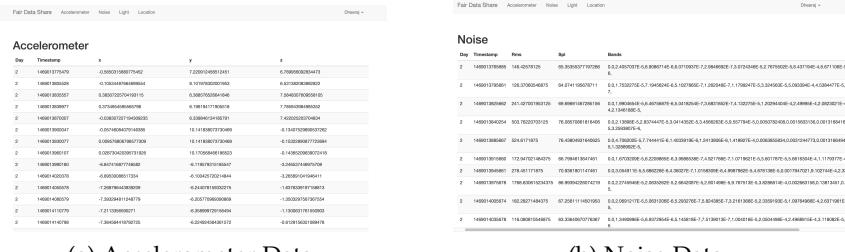


Figure 4.11: User Data

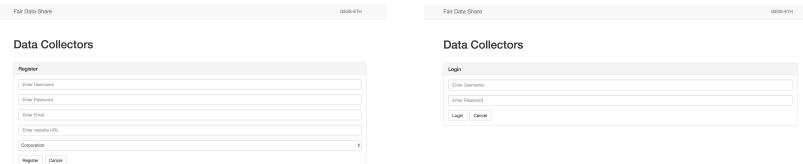


Figure 4.12: Entering the Portal for Data Collectors

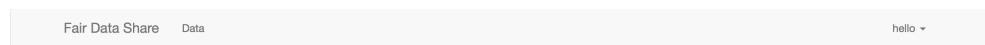


Figure 4.13: Data Collectors Welcome Page

Chapter 5

Explanation of the Mobile Application

5.1 The Building Blocks

explain with diagram what interacts with what

5.2 The Mobile Application

5.2.1 Local Storage

5.2.2 Alarms

Going to the Next Data Sharing Day

Notifications

5.2.3 Privacy and Credit Improvement

5.2.4 Recommendations

how privacy credit is imporved

5.2.5 Recording User Choices

5.2.6 Sensor Data Collection and Summarization

5.2.7 Server Synchronization

5.3 The Server

5.3.1 Kinvey Data Storage

Kinvey is a mobile backend as a service which provides a platform for mobile phones to link applications at a backend cloud storage. This backend has been used to store data and for some business logic implementations.

5. EXPLANATION OF THE MOBILE APPLICATION

Security

Table Store

All the data collected from the user's phones is stored in Kinvey Data Store's collections. Data is segregated into the appropriate collections. Data is stored in the collections is done so in the following form :

1. Data collected in the form of radio buttons on the phone are stored as integers.
2. Data collected as check-box entries on the phone each have a column in the collection with entries zero or one.
3. Data collected as drop down lists on the phone are stored either by integer position in the list or by the entry name in the list itself.

This way of storing data applies to all the collections explained below. Each record of every collection is with a unique user identifier and with the timestamp whenever necessary. This is done to identify data that belongs to the same user across different collections. The timestamp is collected in order to examine temporal relationships. The first collection is the GetUserInformation collection and is used to store all the basic non intrusive user information collected in the entry phase. The collection is shown in figure 4.1b. Next, is the UserResponse collection which is used to store all the responses of the users to data requests in the entry phase and the core phase. The collection is shown in figure 4.1b. Collection AccelerometerStore, LightStore, NoiseStore and LocationStore are used to store the mobile sensor data collected from the user at the end of a bidding day after local summarization on the mobile phone. These collections are shown in figures 4.1b 4.1b 4.1b 4.1b. The StorePoints collection is used to store the privacy and credit metrics obtained at the end of each bidding day for each user and is shown in figure 4.1b.

Bussiness Logic

?? Most of the business logic used in the FairDataShare portal is present on Kinvey. Two things are done here:

1. Finding privacy for a user
2. Summarization

Data collected from the user consists of mobile sensor data with the least amount of summarization. This prevents for repeatedly collecting data from users and saves space and mobile data. Therefore, this mobile sensor data needs to be further summarized before being given to the stakeholder. To do this, we first have to find the most recent privacy setting from the UserResponseCollection. This is done due to the fact users can answer a data

request more than once, hence we need to fish out the latest response for a particular sensor. Once this is done, using the recorded privacy level, we feed this input into another script that performs the summarization which has been explained in chapter 3.

5.3.2 FairDataShare Web Portal

The FairDataShare portal makes use of a server other Kinvey to safely store the usernames, passwords of users and the stakeholders. The database technology used is MongoDB. The username and passwords are both stored in a collection. 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 described in section ???. The webpage was constructed using Html and css. Screenshots of the portal are provided in chapter 4.

Chapter 6

Pre-Survey and Experiment Findings

6.1 Overview of the Pre-Survey Data

basic statistics about the data, plots to explain the data

6.2 Pre-Survey Methodology and Findings

6.3 Overview of the Experiment Data

6.4 Findings from the Experiment Data

Chapter 7

Conclusion

Appendix A

Appendix

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.