

Predicting perceived stress levels with phone application usage

Smartphone Usage as a Predictor for Perceived Stress

Name student: Gary van Koeeverden

Student ANR: 299522

Email: g.c.vankoeeverden@tilburguniversity.edu

THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF
MASTER OF SCIENCE IN DATA SCIENCE & SOCIETY
DEPARTMENT OF COGNITIVE SCIENCE & ARTIFICIAL INTELLIGENCE
SCHOOL OF HUMANITIES AND DIGITAL SCIENCES
TILBURG UNIVERSITY

Thesis committee:

Dr. A. Hendrickson

Dr. M. Jung

Tilburg University School of Humanities and Digital Sciences

Department of Cognitive Science & Artificial Intelligence

Tilburg, The Netherlands

January, 2020

Abstract

Stress has a growing effect on society, predicting stress through smartphone usage seems a cost-effective and convenient method to measure stress. This study investigates the influence of different phone usage features on perceived stress levels and tests if these features can identify perceived stress levels. Five different models are tested, both user-specific and generic models. Three different classification algorithms were developed: Random Forest, Support Vector Machine and k-Nearest Neighbours. The data consisted of two different sets, one dataset that consisted of returned mental health surveys from a group of respondents and the other dataset was the phone usage log data of the same group. This data was merged and the aim of this study was to predict stress from small time frames of maximum two hours of phone usage data upon every returned survey. First an exploratory analysis was done on the different models to test which features have the strongest relation with the target stress levels. Afterwards the five models were tested with the classification algorithms. The classification results indicate that the classification algorithms do not perform better on the user-specific models as predictor for stress than on generic models. The different classification algorithms and models show very dissimilar results and predict in general not better than the baseline.

Contents

Abstract	1
1. Introduction	3
1.1 Context	3
1.2 Problem Statement and Research Questions	5
1.3 Thesis Outline	5
2. Related work.....	6
3. Experimental setup	9
3.1 Dataset Description.....	9
3.1.1 Dataset	9
3.2 Data cleaning	10
3.3 Feature engineering.....	11
3.3.1 Application Categories	12
3.3.2 Feature extraction	12
3.4 Exploratory Data Analysis.....	15
3.4.1 Target feature	15
3.4.2 Application category features.....	16
3.5 Data split	18
3.6 Imbalanced Dataset	19
3.7 Algorithms	19
3.8 Evaluation method.....	22
3.8 Parameter Selection and Tuning.....	23
3.9 Software	24
3.10 Setup	25
3.10.1 Baseline	25
3.10.2 Setup research question 1	25
3.10.2 Setup research question 2.....	26
4. Results	27
4.1 Exploratory analysis	27
4.1 Classification models.....	30
5. Discussion, Limitations & Conclusion	33
5.1 Discussion.....	33
5.2 Limitations and future research.....	34
Bibliography.....	36
Appendix A.....	40
Appendix B.....	45

1. Introduction

This section will introduce the complete study. Section 1.1 provides a clear introductory context and in section 1.2 the problem statement and research questions are stated. Section 1.3 provides an outline for the rest of the thesis.

1.1 Context

These days everyone is extremely busy with life. We have an education, a job, a partner, a family, a social life, a hobby and if there is some time left even sports. It seems as if we have never been busier, there are not enough hours in a day to do all the things we want or need to do. Not only in our personal life we experience a tight schedule, the same is seen at work. Cramming too much activities and work into a short amount of time causes stress (Franke, 2003), which is the number one occupational disease in the Netherlands. In 2018 thirty-six percent of all the Dutch work related sick leave was due to work stress. In this same year over 1.3 million Dutch people showed signs of physical and emotional exhaustion. The total amount of stress related sick days exceeded 11 million, and the total costs of this absence exceeded 2.8 billion euros. This comes down to 8.100 euros annually per employee nationwide (TNO, 2019).

According to Hooftman et al. (2019) there is a trend visible, the amount of people that show stress symptoms has been growing since 2007. Every living creature experiences stress, but also needs it. A limited amount of stress aids the body and helps to cope with different situations. Experiencing too much stress over a longer period causes several health issues. Stress is associated with coronary heart diseases (Rosengren, et al., 2004), immune dysregulation and cancer (Godbout & Glaser, 2006), and reduces brain tissue volumes (Blix, Perski, Berglund, & Savic, 2006).

Early identification and prevention of extensive stress can be very beneficial on a personal level but it might also be interesting for organisations too. Kompier & Cooper (2003) showed that reducing the amount of stress in the workplace is a means to reduce employee costs. Furthermore it also improves the working environment within companies. Stress reduction can be very profitable on multiple levels, and a tool to measure stress might be closer than thought.

Smartphones have become integrated into everyday life, one in three Dutch people cannot even go to the toilet without bringing the smartphone along. With millennials this number is even higher: two out of three millennials bring their smartphone along to the toilet. Since the smartphone is always present in daily life, it seems to be an excellent tool to monitor daily

activities. In 2019 Dutch people daily spent on average two hours and fifteen minutes on their phone, this comes down to 34 full days of smartphone usage per year (Simyo, DirectResearch, 2019). The smartphone has become our small personal assistant, and collects lots of data. In this study it is tested if phone usage behaviour can predict stress. The goal of this study is to explore what phone usage features have an influence on the perceived stress levels of phone users, and tries to predict perceived stress levels based strictly on phone application usage in short time windows of maximal two hours. Another aim of this research is to test what kind of model the best predictor of stress is, a user-centric model or a generic model. This is done with the use of three different classification algorithms: Support Vector Machines, Random Forest and k-Nearest Neighbours.

In the past there already has been done research on phone usage as a predictor of stress. Multiple studies have focused solely on phone application usage (Ahn, Wijaya, & Esmero, 2014; Ferdous, Osmani, & Mayora, 2015), some included Bluetooth interactions (Bogomolov, Lepri, Ferron, Pianesi, & Pentland, 2014), other also included call logs, Wifi and 3G connectivity and locational data (Sarker, Kayes, & Watters, 2019) (Kang, Seo, & Hong, 2011). There has also been research on stress identification with the use of smartphones and wearables (Zenonos, et al., 2016). These wearables collected data through physiological and movement sensors. These are all very pervasive methods that ask a lot of the privacy of the phone users.

From a societal point of view, predicting stress through phone usage is a very cost-effective and uncomplicated approach. The smartphone penetration in The Netherlands is 93% (Deloitte, 2019). More than 9 out of 10 people have a smartphone. Since the costs of stress have been rising over the years, an easy tool to measure and identify stress might help to fight these expenses. This cost reduction is very interesting for both the government and the industry. Also on a personal level there are benefits that derive from stress identification. This study can help identify the important phone usage features and might aid for future development of stress identification applications. From a scientific point of view, this study has an approach which has not been used before in similar research. The focus lies on small time windows and purely on the phone application data. The identification of the phone usage features that have a contribution to the performance of classification model helps narrowing down the scope for future research on stress detection through phone usage.

1.2 Problem Statement and Research Questions

Problem statement:

To what extent can perceived stress levels be predicted by phone application usage?

Sub questions:

Q1. Is there a distinction between the factors that influence the perceived stress levels for the user-specific models and the generic models?

Q2. What model predicts best the perceived stress levels of smartphone users?

Five different models will be developed to answer the research questions. Three of these models will be user-specific models and be based on the data of three randomly chosen respondents. Another model is based on the top30 respondents with the most completed survey responses and the last model is based on all respondents¹. An exploratory data analysis will be done to answer research question 1. For research question 2 different machine learning classification models will be used to determine the optimal model for predicting stress.

The results are very differing for each of the models and the classification algorithms. The correlation between the different phone usage features and the perceived stress levels are moderate to strong for the user-generic models. While for the generic models there are no features that are strong linearly correlated. The balanced accuracy scores for the performance on the test sets for the user-specific models vary between 96.7% and 46.4%. While the generic models show less variance, these balanced accuracy scores vary between 50.6% and 61.2%.

1.3 Thesis Outline

The remaining of this thesis research paper is structured as follows: Section 2 reviews previous studies on stress detection and the relation of smartphone usage with stress. Section 3 describes the datasets that are used, feature selection and model development. The experimental results are shown in section 4 and discussed and concluded in section 5.

¹ These models are developed on a subset of the data, the respondents are also randomly selected from this subset. This is clarified in section3.3

2. Related work

This section will provide an overview of previous research on stress detection and phone usage.

2.1 Stress detection

A tight schedule, a car accident and a stock market crash all seem very different at first sight. But there is at least one thing that they have in common: stress. Even though the understanding of stress in these three situations is completely different, we still experience the same thing, stress. Since it is such a catchall term, it makes it very hard to define. Stress can be approached from different angles, and has already been addressed within a broad field of different sciences. The fields of biology, psychology and sociology all have their own definition of stress. (Epel, et al., 2018) This is mainly due to the fact that stress can be measured on multiple levels and through different ways.

There are several physiological measurement methods to analyse stress. Measuring cortisol levels (Kirschbaum, 1993), brain tissue volumes (Blix, Perski, Berglund, & Savic, 2006) and blood pressure or heart rates (Lebepe, Niezen, Hancke, & Ramotsoela, 2016) are ways to measure stress. All these are common practices and are highly reliable ways to measure stress but they do require the use of expensive equipment and technologies. Next to that, these technologies most often need professional guidance to be operated.

Wearables

Zenonos, et al. (2016) developed a framework in which they used smartphones and wearables to recognize mood states at the workplace. Four participants were monitored for a total of 11 workdays, between 09:00 – 17:00, with the use of a wristband and chest sensor. These devices were embedded with several physiological sensors that registered heartrate, pulse rate, body temperature and acceleration. Every two hours the respondents were asked to fill in a survey on their smartphone regarding their overall mood in that time frame. The researchers achieved a 62% accuracy on a generalized classification model and 69% accuracy on personalized classification models that predicted stress.

Smartphone usage

Ferdous, Osmani, & Mayora (2015) focused on stress recognition on the workforce, and monitored the smartphones of 22 participants for over 6 weeks. In their study they collected the phone application usage of the respondents and prompted a survey three times a day. In this survey the respondents were asked what their perceived stress level was. The researchers categorized the used applications into 5 different categories: “Social Networking Service Applications”, “Entertainment applications”, “Utility applications”, “Browser applications” and “Gaming Applications”. The total amount of unique applications that were recorded during this study was 128. Participants used on average 12 different applications, with a standard deviation of 6.45. In this study the researchers used a metrics-based approach, where they used count data for classification. The features that were extracted from the data are the amount of unique applications used, time spent in every application category and the amount of times the applications from the different categories were used. The authors reached a 54% accuracy with a generic support vector machine (SVM) classification model that used the data of all the respondents combined. With the use of cross-validation and a user-centric model the researchers reached an average accuracy of 75%. The researchers showed with their work that the respondents use their smartphone disparate. In this study the approach will be sort-like, and the applications will be categorised into slightly more categories. This is due to the fact that applications are more widely used these days.

Bogomolov, Lepri, Ferron, Pianesi & Pentland (2014) have reported that they are able to recognize daily stress reliably based on smartphone data. In their research they monitored 117 participants for more than 8 weeks. During this time, all the respondents received an android smartphone with specialized software that registered call and sms data and social proximity data. With the use of Bluetooth the software scanned for near-by devices. The authors also collected daily surveys regarding the respondents personality (“Big Five” personality traits) and the experienced daily stress. The researchers argued that the weather conditions have impact on daily mood and included the daily weather conditions in the data. The authors approached the stress recognition task as a binary classification problem and reached an accuracy of 72.39% with a generic model that combined all the user’s data. The transformation to a binary classification task will be adopted into this research.

Bauer & Lukowicz (2012) have shown that students show different mobile phone usage patterns when stressed. Their study focused on locational, call and SMS behaviour and left all other features out. Authors of another study managed to develop a similar model that reached an F-score of 74.2%. The participants in this study were asked to fill in a survey once a day and their smartphones were logged. Also the data of the sensors was logged, which gave insights into illuminance, acceleration and orientation of the phone, e.g., how parallel the smartphone was to the ground or how bright the room was. Two other studies on smartphone usage patterns show that stress can be better predicted for an individual than for a group of people (Bauer & Lukowicz, 2012) (Muaremi, Arnrich, & Tröster, 2013). The individual usage patterns of the smartphone users show a lot of variation. This suggests that it might be arduous to generalise phone usage to a larger population, and a user-centric model might suffice more than a generic model.

3. Experimental setup

This section provides information about the raw datasets that are used for this study. Furthermore this section will also provide an exploratory data analysis, giving a description on the features of the dataset. Thereafter the features are selected for the Machine Learning Experiments, along with the training of the classifier procedures and the performance criteria.

3.1 Dataset Description

For this study, two anonymized datasets (Hendrickson, Abeele, & Aalbers, (under review)) in Comma Separated Value (CSV) file extension have been provided. The datasets contain data about a group of people that volunteered in a research on phone usage. For a period of approximately four weeks the smartphones of these participants were logged with use of the MobileDNA application. This application logs phone usage and metadata. Next to that the respondents were regularly, with 4 times per day, asked to fill in a survey regarding their current activities, social interactions and different moods. To receive the surveys, the respondents had to install another application, Ethica. The phone log data is collected in the Phone Usage Dataset and the survey responses are collected in the Mood Survey Dataset. For the sake of clarity, in the rest of this thesis the phone usage data will be referred to as phone data and the mood usage data as mood data.

3.1.1 Dataset

The phone dataset consists of 124 unique user ID numbers. These user IDs used the applications for 586.792 times in 149.889 phone sessions between 21-02-19 and 26-3-19. The Phone Dataset consists of 586.792 rows with values that are separated by a comma. One row of data represents the log of single application that is used by a user.

Feature	Description	Type
Application_ID	Unique name of application that is used	Character string
Battery_level	Battery level in percentage (%)	Ratio
End_time	Time that application was closed	Timestamp
End_time_MilliS	Conversion of End_Time into milliseconds	Characterstring
Notification	Whether a notification initiated the application	Binary
Session_ID	ID number for every unique phone usage session	Character string
Start_time	Time that application was opened	Timestamp
Start_time_MilliS	Conversion of Start_Time into milliseconds	Character string
User_ID	Unique ID number of respondent	Character string

The following string is an example of a row a values, extracted from the dataset. Nine variables, separated by a comma, combined in one string. Table 3.2 gives an overview of the different features in the dataset.

[com.ethica.logger, 96, 2019-03-12T12:09:11.390, 1552388951390, False, 1552388935, 2019-03-12T12:09:02.630, 1552388942630, 12345]

The mood dataset consists of 16.016 entries from 149 unique user id numbers, collected between 04-06-18 and 14-05-19. The id numbers from this dataset correspond to the id numbers of the phone data set. This file consists of 16.016 rows with values that are separated by a comma. One row of data represents a survey that is returned by a user.

Table 3.3 below gives a concise overview of the features in the mood dataset. The following string is an example of a row of values, extracted from the dataset. Thirty-five variables, combined in one string. *[12345,2019-02-22 10:39:41 CET, 2019-02-22 10:41:22 CET, 1, 4, 3, 3, 3, 4, 1, 2,5,2,3,1,0, Working ,Together with coworkers ,3,0,0,1,0,0,0,0,0,0,0,0,0,0,1,1]*

Table 3.3	35 Features Mood Dataset	
Feature	Description	Type
User id	Unique ID number of respondent	Categorical
Sent_time	Time when survey was sent to respondent	Timestamp
Resp_time	Time when survey was returned by respondent	Timestamp
Duration	Duration – time used to fill in survey in minutes	Ratio
12 mood variables	12 different mood categorical variables	Categorical
Activity	Recent activity of respondent, categorised in 9 different options	Categorical
Social_enjoy	Whether the respondent liked current social setting	Categorical
6 social setting variables	The social feature transformed into dummy variables	Binary
9 activity variables	The activity feature transformed into dummy variables	Binary
Time_window	What part of day the survey was sent, day divided in 4 windows (1-4)	Categorical

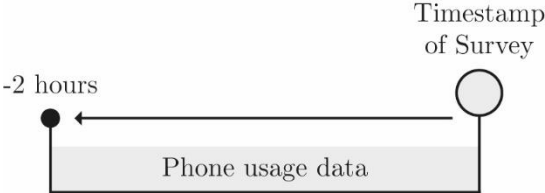
3.2 Data cleaning

The datasets both contain an amount of erroneous duplicate data. This data is removed from these datasets before the data was further pre-processed. This process contained of deleting duplicates, erroneous user ID's, expired, double, blocked and cancelled surveys. This process has decreased the total number of surveys to 6.323. The process can be found in Appendix B.

3.3 Feature engineering

The respondents were very inconsistent with returning a survey. As a result of the inconsistency the respondents frequently returned multiple completed surveys per day but also had several days of not returning a single survey. Therefor the chosen approach in this study is one that focusses on the time frame of two hours upon the surveys. For this study all the returned filled-in surveys are used as a single data measure point. A time frame of two hours was created with the survey response time as the endpoint. All the phone usage data in the two hours until the endpoint was collected. Figure shows a visualisation of this method. The surveys were approximately sent between 08:00 am and 11:00 pm. This is a total time window of 15 hours wherein the respondents can return 4 different surveys. The surveys were not sent at the exact same point of time daily, also the respondents had two hours to complete the survey. As a result of this some respondents returned multiple surveys within a very small time frame. Overlapping data is not beneficial for data analysis with relatively few observations (Harri & Brorsen, 2009). For this reason all the returned surveys from a single respondent within two hours of another survey are omitted (n=584). So there are no overlapping time frames. On the opposite, a number of respondents did not use their phone in the two hours upon replying the survey, these surveys (n=555) are also excluded since there is no corresponding phone data.

Fig 3.1 Visualisation of the process of collecting data



3.3.1 Application Categories

As a first step all the applications² from the phone usage dataset are categorized into 10 different categories. These categories are based on the categories in which the applications are found in the Google App Store. Application usage is based on popularity, some application categories are not used as often as others. Therefore some categories are combined, as seen in earlier research (Ferdous, Osmani, & Mayora, 2015). Table 3.4 gives an overview of the different categories.

Category	Feature name	Description
Entertainment	entertainment	All entertainment apps
Finance	finance	All financial and governmental apps
Games	games	All gaming apps
Lifestyle	lifestyle	All news, lifestyle and travel apps
Online shopping	onlineshopping	All online shopping apps
Process	process	All background processes
Productivity	productivity	All productivity apps
Social	social	All social media and communication ³ apps
Utility	utility	All utility apps
Wrong	wrong	Undefined/ unknown applications

3.3.2 Feature extraction

For this study there are 28 features extracted from the data. This study uses a metrics based approach, so these features are mostly based on count data. The features can be divided into three groups: application category variables, time variables and other count variables. The features are described on the following two pages.

Features

1. Total time spent in application from category.

(Ente_time, Fina_time, Game_times, Life_time, Shop_time, Proc_time, Prod_time, Soci_time, Util_time)

These 9 features measure the total time spent in the applications from that category. The nine categories are Entertainment, Finance, Games, Lifestyle, Online shopping, Process, Productivity, Social Media and Utility. The total time is measured in seconds.

² The Ethica application data is included in the data and placed in the utility category.

³ The applications from these categories are combined since there is a thin dividing line that separates the categories

2. Total count of times application used from specific category

(Ente_count, Fina_count, Game_count, Life_count, Shop_count, Proc_count, Prod_count, Soci_count, Util_count)

These 9 features measure the total amount of times the applications from that category are used. The nine categories are Entertainment, Finance, Games, Lifestyle, Online shopping, Process, Productivity, Social Media and Utility.

3. Total amount of notification initiated sessions (noti_count)

This feature is the count of the sessions that are initiated by a notification.

4. Total time on phone (total_time_phone)

This feature is created by measuring the total amount of time that the phone has been in used. The total time is measured in seconds.

5. Session count (sess_count)

This feature is the total count of sessions in the time frame.

6. Minimum length session (min_len_sess)

This feature is the length in seconds of the shortest session in the time frame.

7. Maximum length session (max_len_sess)

This feature is the length in seconds of the longest session in the time frame.

8. Mean length session (avg_time_sess)

The mean time of a session in seconds per time frame.

9. Unique applications (uni_app_count)

Amount of unique applications that are used in a time frame.

10 Weekend (if_weekend)

A binary feature that indicates if the survey was completed during the weekend.

11 Hour of day (xhour)

Time of day in hours transformed with a trigonometric function. This is a measure to address

the cyclicity of time. The result of this transformation is that for example a value of 23 lies closer to 2 than to 20, which is the case with time. The function is formulated as:

$$xhour = \sin\left(2\pi \frac{h}{24}\right)$$

Where sin is sine, π is pi and h is the hour of the day in the 24 hours notation.

12. Hour of day (yhour)

Time of day in hours transformed with a trigonometric function. This is a measure to address the cyclicity of time. The function is formulated as:

$$yhour = \cos\left(2\pi \frac{h}{24}\right)$$

Where cos is cosine, π is pi and h is the hour of the day in the 24 hours notation.

3.4 Exploratory Data Analysis

3.4.1 Target feature

Binary classification task

Table Classes count past and present situation			
Multiclass classification		Binary Classification	
Class 0	1013	Class 0	2728
Class 1	933	Class 1	1074
Class 2	782		
Class 3	607		
Class 4	374		
Class 5	93		

Every person experiences and perceives stress differently. Two people can experience the same amount of stress but perceive it differently. This is also seen in the distribution of the survey answers of the respondents. Some respondents perceive stress level 1 ‘very slightly’ as the baseline, where others perceive stress level 0 ‘not stressed’ as their baseline.

There is only a limited amount of survey data per respondent, next to that there is a high class imbalance for a considerable amount of the respondents. Table shows the distribution of the classes. Seven respondents only perceived 1 level of stress and four respondents perceived 2 levels of stress.

There are 84 respondents that have chosen a perceived stress level category only one time in a survey, and 46 respondents that have chosen a perceived stress level two times. It is impossible for a classification model to classify with classes that have only one entry. In the case of very small classes merging is seen as a desirable solution. (Fukunaga, 2013) Classes that are similar can be merged. Therefore the six classes have been merged into two larger classes. Class 0 is ‘not stressed’, which is the combination of stress levels 0-1. Class 1 is ‘stressed’, with stress level classes 2 – 6 combined. Figure 3.4.1 displays the previous class distribution and figure 3.4.2 shows the current distribution.

Figure 3.4.1 Old class distribution

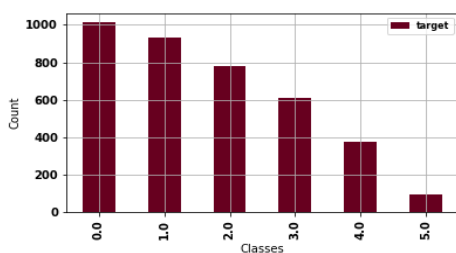
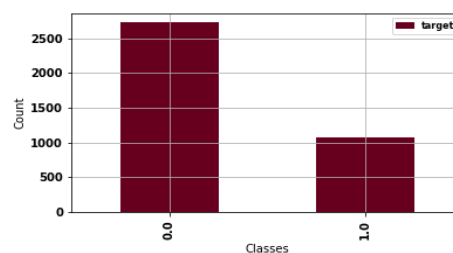


Figure.3.4.2 New class distribution



3.4.2 Application category features

The applications are placed in nine different categories and for every category there are two features that have been extracted from the data. This is the total time spent in the application and the count of the times that applications from that category have been used. The descriptive statistics of the time features are displayed in table 3.4.2.1.

Feature	Count	Mean	Std	Min	25%	50%	75%	max
Ente_time	3802.00	<u>277.96</u>	778.13	0.00	0.00	0.00	66.00	6583.00
Fina_time	3802.00	5.39	31.94	0.00	0.00	0.00	0.00	997.00
Game_time	3802.00	47.45	273.53	0.00	0.00	0.00	0.00	4182.00
Life_time	3802.00	66.41	310.62	0.00	0.00	0.00	7.00	6334.00
Prod_time	3802.00	40.86	191.82	0.00	0.00	0.00	18.00	4976.00
Proc_time	3802.00	8.78	35.68	0.00	0.00	1.00	5.00	825.00
Shop_time	3802.00	5.93	56.65	0.00	0.00	0.00	0.00	1186.00
Soci_time	3802.00	<u>744.69</u>	921.18	0.00	116.00	420.00	1006.75	6749.00
Util_time	3802.00	<u>249.40</u>	462.00	0.00	52.25	97.00	228.00	5943.00

The descriptive statistics from table 3.4.2.1 suggest that the data is skewed. A lot of zeroes are registered. Which seems normal since the average amount of applications installed on a smartphone lies between 60 and 90 (Annie, 2017). As can be seen from table 3.4.2.2, is the mean of uni_app_count 7.18. This indicates that that on average, around 10 percent of all applications on a phone are used in the time frames. As a result, a lot of zeroes are counted.

Feature	Count	Mean	Std	Min	25%	50%	75%	max
Noti_count	3802.00	3.25	3.88	0.00	1.00	2.00	4.00	44.00
Uni_app_count	3802.00	7.18	3.72	1.00	4.00	7.00	10.00	25.00

As can be seen from table 3.4.2.1, the mean time spent in a category differs a lot between the categories. The mean time spent in applications from the Entertainment, Social Media and Utility categories is significantly higher than the other categories. Not only the mean time shows a lot of variation, also the total time spent in the different categories varies substantially.

Figure 3.4.2.1 Popularity category total time count Figure 3.4.2.2 Total count application used

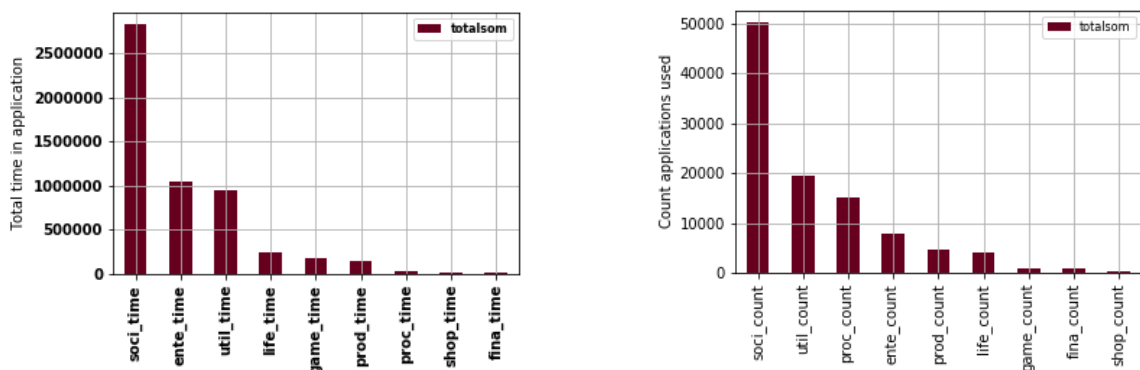


Figure 3.4.2.1 and 3.4.2.2 show the popularity of the different application categories. The social media category is the most popular one, and applications from this category are more often and longer used than the other categories combined. Social media applications are used often but for very short amount of time. This is in line with the results from the research by Oulasvirta, Rattenbury, Ma, & Raita (2012) which claims that people like to check social media regularly to receive a little informational reward.

Figure 3.4.2.3 Mean time spent per use

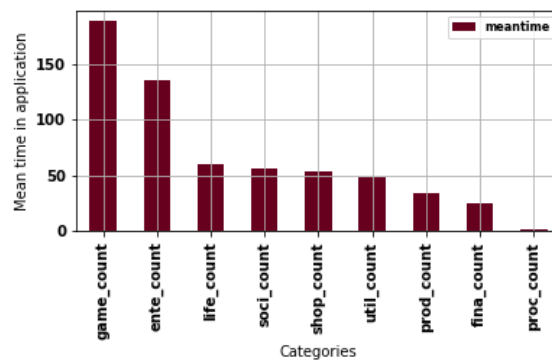


Table 3.4.2.3 displays the descriptive statistics from the features that are related to the sessions.

Feature	Count	Mean	Std	Min	25%	50%	75%	max
Min_len_sess	3802.00	93.35	533.89	0.00	1.00	5.00	27.00	7144.00
Max_len_sess	3802.00	782.81	1044.12	0.00	146.00	391.00	971.00	7144.00
Avg_time_sess	3802.00	<u>285.71</u>	624.87	0.00	55.00	107.00	246.12	7144.00
Total_time_phone	3802.00	<u>1446.85</u>	1444.31	0.00	317.00	963.50	2117.50	7144.00
Sess_count	3802.00	<u>8.72</u>	7.32	1.00	3.00	7.00	12.00	68.00

The respondents used their phone on average 8.72 times per time frame. The mean total time that a phone was used per time was 1446.31 seconds, this comes down to 24.1 minutes. After unlocking the phone, the respondents spent on average 4 minutes and 45 seconds per session.

Figure 3.4.2.4 Application category usage per model

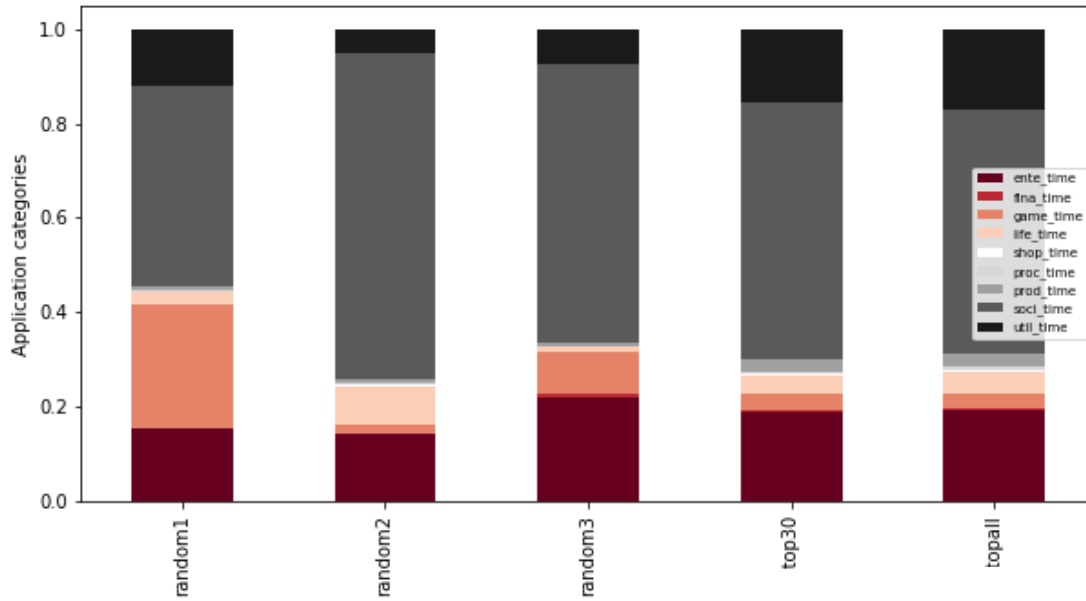


Figure 3.4.2.4 is a visualisation of the relative spent time in each application categories, based on the specific data for the earlier described models which are later used for the classification task. The distribution for the Top30 and TopAll models are fairly similar, but the individual models Random1, Random2 and Random3 show less similarity.

3.5 Data split

The performance of the different models needs to be tested. To test the performance, the researcher needs different data sets for training and testing. The creation of these datasets is done by splitting. The data for all the different models will be split into two sets, a training set (80%) and a test set (20%). This data is first randomly shuffled with the use of the shuffle function from the sklearn library, after this the data is split with the stratified sampling technique from that same library. Stratification is used to make sure that both the training and test set contain entries from both classes (Köhl, Magnussen, & Marchetti, 2006).

3.6 Imbalanced Dataset

As discussed in the previous sections, the dataset is fairly imbalanced. Class 1, the stressed class, has the lowest number of entries. There are multiple approaches to solve this inequality. Beyan & Fisher (2015) suggest that there are four approaches: on algorithmic level, data level, cost-sensitive methods and ensembles of classifiers. In this study the chosen approach is on data level. On this level there are mainly two options: undersampling and oversampling. With undersampling a subset of randomly selected entries from the majority class is created, this subset matches the size of the minority class (Maimon & Rokach, 2005). This produces a more homogeneous dataset. The risk of undersampling, is that it may omit entries that contain useful data, and this can have a negative influence on the performance of the classification models. With oversampling, random entries from the minority class are duplicated and included in the majority class. (Tan, Steinbach, & Kumar, 2014), a downside of this is that the classification model will overfit on the training data. To address this problem, another technique has been developed: SMOTE (Chawla, Bowyer, Hall, & Kegelmeyer, 2002). This is an abbreviation of Synthetic Minority Over-Sampling Technique. As the name suggests, this method produces synthetic entries by interpolating between multiple entries from the minority class. This technique produces samples that are not exact duplicates but approximate combined copies with extra noise added. SMOTE is a powerful method that prevents overfitting, therefore this technique will be applied on the training data in this study.

3.7 Algorithms

Classification models

Wolpert, Macready & others (1997) state in their work, the ‘No Free Lunch’ (NFL) theorem, that there is no optimal machine learning model for every problem. Every problem is unique and has its own assumptions, and therefore the model that works best is highly dependable on the problem itself. Therefore three different classification algorithms will be used in this study: Random Forest, Support Vector Machines and K-Nearest Neighbours.

Random Forest

The Random Forest algorithm is based on a collection of decision trees. A decision tree (DT) is one of the most basic classification algorithms that is commonly used. A decision tree is a collection of decision nodes that are connected with branches. Starting from the root node, features are tested at and binary split on every decision node and branches derive from these nodes. One branch for each possible outcome, at the end of this branch there is a new decision node unless this is a leaf node. There are various measures to split a node, in this study the Gini Index is used. The Gini Index is formulated as

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$$

\hat{p}_{mk} is the proportion of training observations in the m th region that are from the k th class. The Gini index is also seen as the measure of node ‘purity’. When the Gini index is a small number, it implies that the node contains mostly observations from one specific class. (James, Witten, Hastie, & R., 2017)

Random forests are a combination of multiple decision trees. Every single tree in ‘a forest’ depends on the values of a random vector that is independently sampled. All these trees are sampled with the same distribution. The number of trees used in a random forest has influence on the generalization error, but converges to a limit when the number of trees becomes substantial (Breiman, 2001). The difference between a decision tree and a random forest is that a decision tree is hierarchical and chooses as the first input variable the one that explains the most variance. With a Random Forest the decision trees are fed random input variables, this is a measure to reduce the variance. (James, Witten, Hastie, & R., 2017)

K- Nearest Neighbour

K-Nearest Neighbour (k-NN) is a so-called non-parametric learning algorithm that can be used for regression and classification problems (Altman, 1992). This is an instance-based algorithm which has the training data saved in it. When a new unseen data entry is fed to this classification model, it will match with the k nearest neighbours in the feature space. A majority vote follows by its k neighbours. The entry is then classified as the class which is the most common amid the k neighbours. The value of k , which is a positive integer, has influence on the decision boundary. A low value of k can lead to overfitting, this is because the decision boundary is more complex. A high value of k can lead to underfitting, since the decision

boundary is simpler. (Larose & Larose, 2014) There are several distance measures to measure similarity. In this study the Euclidean distance is used, which is formulated as

$$d_{Euclidean}(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Where $\mathbf{P} = p_1, p_2, p_3, \dots, p_m$ and $\mathbf{q} = q_1, q_2, q_3, \dots, q_m$ represent the m feature values of two observations.

Support Vector Machine

The Support Vector Machine (SVM) is a robust algorithm that can be used for regression and classification. A SVM is trained to identify the maximal margin hyperplane in both non-linear and linearly dividable data (Tan, Steinbach, & Kumar, 2014). This hyperplane divides the different classes and the decision boundary is based on the training data. It is also called a soft margin classifier, which implies that the algorithm allows entries to be on the wrong side of the hyperplane. These wrong entries correspond with the training data that has been misclassified by the SVM. The observations from the training data are used as the support vectors, and help minimizing the classification error. SVM avoids the curse of dimensionality and performs successfully with high dimensional data. (James, Witten, Hastie, & R., 2017)

Pearson Correlation Coefficient

To measure the relationship between features, the Pearson Correlation coefficient(r) will be used. It is a measure to calculate the correlation between features. The value of r can vary between -1, which is a perfect negative relation, and 1, which is a perfect positive relation. Values of r near 0 indicate a very poor relation. The formula is defined as:

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Where \bar{X} and \bar{Y} are the sample means for the two features. In earlier research (Evans, 1996) there has been developed a guideline for the strength of correlations. These are: a. very weak: ($r=0.00-0.019$), b. weak: ($r=0.20-0.39$), c. moderate: ($r=0.40-0.59$), d. strong: ($r=0.60-0.79$) and e. very strong: ($r=0.80-1.0$). These guidelines will be used in the study.

3.8 Evaluation method

Accuracy

A much used evaluation metric for classification is accuracy. Accuracy is the proportion of correctly classified true predictions among the total number of predictions.

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{True Negatives} + \text{True negatives} + \text{False Positives} + \text{False Negatives}}$$

The definitions used in the above formula are defined below at the confusion matrix section at table 3.8.1 the bottom of the page.

Balanced Accuracy

There is a risk with using only accuracy as an evaluation metric, this risk is the so-called accuracy paradox. When data is skewed, a classification algorithm can achieve a high accuracy but still have a poor performance. Due to the imbalance of the dataset also the balanced accuracy metrics is included in this study. This metric is defined as:

$$\text{Balanced Accuracy} = \frac{\text{True Positive Rate} + \text{True Negative Rate}}{2}$$

Where the true positive rate is: $TPR = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

The true negative rate is: $TNR = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$

Confusion Matrix

A clear practice to visualise the numbers used to calculate the accuracy and balanced accuracy is the so-called error/ confusion matrix. This matrix gives a clear overview of how the algorithms have classified the data.

Table 3.8.1 – Confusion Matrix		
	Class 0 prediction	Class 1 prediction
Class 0 Actual	True Negatives	False Positives
Class 1 Actual	False Negatives	True Positives

3.8 Parameter Selection and Tuning

The three different classification algorithms have several parameters which can be tuned to improve the performance. In this study the focus will lie on tuning the parameters that can be found in table 3.8.1. Every algorithm has different parameters that can be tuned. With the use of GridSearchCV, which is part of the sklearn library, the optimal parameters can be found. The function tunes the parameters of the algorithms systematically with pre-set values. The performance metrics that it seeks to optimise is accuracy. The GridSearchCV function will split the training data with the use of cross validation into five different partitions. These partitions are then alternately used as training and validation set. Having five partitions means that the GridSearchCV function will perform five different tests, where every partition will function one time as a validation set. Cross Validation prevents against over-fitting and selection bias (Cawley & Talbot, 2010). The GridSearchCV functions tunes and tests the parameters for every partition and selects the optimal ones.

Classification model	Parameters
Random Forest Classifier (RF)	Depth Number of trees Minimum samples split Minimum samples leaf
Support Vector Machine Classifier (SVM)	C: regularization parameter. Penalty error term. ϵ : epsilon Kernel: Radial Basis function (RBF)
K-Nearest Neighbours Classifier (kNN)	K : number of neighbors

The random forest has several parameters, one is depth, which is the maximum depth that a tree is allowed to have. Setting a maximum prevents the tree from overfitting on the training data (Shu, 2020). A larger number of trees improved the performance, but there is a threshold on which the performance gain is not significant anymore. From this moment an extra tree will only cost extra computational power (Oshiro, Perez, & Baranauskas, 2012). Minimum samples split and minimum samples leaf are parameters that prevent against overfitting (Brefeld, 2019). The support vector machine has parameter C, which is a regularization parameter. This parameters limits the influence that the support vectors have on the algorithm (Witten, Frank, Hall, & Pal, 2016). The kernel in this study is set on RBF since it is proven to be the most effective kernel for splitting non linearly separable data (Howlett & Jain, 2001). The epsilon is

the boundary around the classification function, all the erroneous data points that fall within this boundary are ignored, this prevents against overfitting (Hamel, 2011). The value k from the k -NN algorithm is already discussed in section. It is the amount of neighbours that the algorithms should take in consideration to predict the class.

3.9 Software

This study used Python (version 3.6.9) to build models and perform analyses. Preprocessing is done with the use of Jupyter Notebooks that are hosted by the Google Colaboratory servers. Within Python, a variety of packages has been used, these are outlined below. The ‘numpy’ and ‘pandas’ packages are used for data preprocessing. With the use of ‘sklearn’ package the various classification models are tuned and applied. For the data visualizations the ‘Matplotlib’ and ‘Seaborn’ packages have been employed.

Package name	Version	Source
Numpy	1.17.4	(Oliphant, 2006)
Pandas	0.25.3	(McKinney, Wes, & others, 2010)
Scikit-learn/ sklearn	0.21.3	(Pedregosa, et al., 2011)
Matplotlib	3.1.2	(Hunter, 2007)
Seaborn	0.9.0	(Waskom, et al., 2018)

3.10 Setup

There are different datasets chosen for the classification. These datasets are Random1, Random2, Random3, Top30 and TopAll, these sets will be referred to as the five different models. These five different models are used to test whether a generic classification model outperforms models that are user-specific. Respondents random 1, 2 and 3 are randomly chosen respondents from the top30 dataset. These respondent datasets all have at least 80 different data entries. The top30 is the dataset with the entries from the thirty respondents that have returned a completed survey the most often. Top30 is a dataset with 2.447 data entries. TopAll is the complete dataset with entries from all the respondents (n=64) in the sample dataset. Topall consists of 3.802 data entries.

3.10.1 Baseline

The baselines for the different models are based on the size of the majority class of these models. The classifiers will be evaluated with the accuracy and the balanced accuracy metrics. The accuracy baseline will be set to the proportional size of the majority class of the different models. All the different baselines for every model can be found in table 3.10.1. Since the dataset was split with the stratified sampling technique, both the training and test set are given the same baseline. Class 0 is for all the models the majority class

Respondents	Train	Test
Random1	81.7%	81.7%
Random2	63.8%	63.8%
Random3	76.3%	76.3%
Top30	69.6%	69.6%
All Respondents	71.8%	71.8%

3.10.2 Setup research question 1

Research question 1 is stated as follows: is there a distinction between the factors that influence the perceived stress levels for the user-specific models and the generic models? To test the influence of these features on the stress levels, this research has made use of the Pearson correlation coefficient. With the use of the programming language Python all the correlation coefficients of all the extracted features, which are described in section, are displayed in the correlation matrices.

3.10.2 Setup research question 2

Research question 1 is stated as follows: what model predicts best the perceived stress levels of smartphone users? Three different classification algorithms have been trained to predict perceived stress from the respondents. The six stress levels have been transformed into two classes, not stressed and stressed. With the use of the Random Forest, k-Nearest Neighbours and Support Vector Machine classifier algorithms the different models are tested. The data was split in 80% training data and 20% test data. The parameters of the various algorithms have been tuned with the use of GridSearchCV, which cross validated the training into five folds. The parameters were tested on the test dataset, and these performance scores are measured in accuracy and balanced accuracy. The scores are compared to the baseline score, which is related to the size of the majority class.

4. Results

This section will provide the results from the various analyses. Section 4.1 will give the results from the exploratory analysis and in section 4.2 the results from the classification tasks will be provided.

4.1 Exploratory analysis

The correlations between the features for respondent random1 are shown in Table 4.1.1. This respondent has not installed or used any finance or online shopping applications during the test period. Therefore the fina_time, shop_time, fina_count and shop_count features are omitted. As can be seen from the table has the target feature the strongest correlation the features ente_time, max_len_sess and avg_time_sess, respectively, with moderate correlations of 0.4188, 0.5201 and 0.4916. A noteworthy addition, both the avg_time_sess and total_time_phone features are very strongly related with each other and other features. The features and correlations of respondent random1 will be further discussed in section 5.

Table 4.1.1 - Random1, Correlation matrix 10 most influential features

	soci_time	util_time	prod_count	total_time_phone	sess_count	max_len_sess	avg_time_sess	if_weekend	yhour	target
ente_time	-0.1418	-0.0431	-0.0793	0.6058	-0.1830	0.7973	<u>0.8289</u>	0.2179	-0.1549	<u>-0.4188</u>
soci_time		0.0351	0.1622	0.3654	0.3616	0.1906	0.0958	-0.0341	-0.0582	-0.1472
util_time			0.0846	0.2825	-0.0269	0.2334	0.2641	0.1497	-0.0590	-0.2230
prod_count				-0.0039	0.4255	-0.1454	-0.1334	-0.1956	0.1459	0.1777
total_time_phone					0.0702	<u>0.8817</u>	<u>0.7809</u>	0.2595	-0.0338	-0.3606
sess_count						-0.2031	-0.2922	-0.1296	0.3763	0.2638
max_len_sess							<u>0.8986</u>	0.2882	-0.1930	<u>-0.5201</u>
avg_time_sess								0.2748	-0.2376	<u>-0.4916</u>
if_weekend									-0.0008	-0.2199
xhour										0.1988

Table 4.1.2 shows the ten strongest correlations between the different features for respondent random2. The sess_count, xhour and proc_count are the features that have the strongest correlation with the target feature. The correlation coefficient for these features are 0.3037, 0.2420 and 0.1923 respectively, which are weak relations. Proc_count is one of the features that is stronger correlated with the target variable, but is also highly correlated with soci_count and sess_count. Also, fina_time is very strongly correlated with fina_time. A very strong relation is also seen between soci_count and sess_count ($r = 0.822$). This will be further discussed in section 5.

	fin_a_ time	life_ time	fin_a_ count	shop_ count	proc_ count	soci_ count	sess_ count	xhour	yhour	target
ente_time	-0.0617	-0.0558	-0.0852	-0.0412	-0.0120	-0.0754	-0.1808	-0.0556	-0.0845	-0.1223
fin_a_time		-0.0157	<u>0.8218</u>	0.1512	0.0369	0.1328	0.2278	0.0790	0.0007	0.1348
life_time			-0.0003	0.0074	0.0253	-0.0910	-0.0656	-0.1113	-0.0584	0.1638
fin_a_count				0.1437	0.1380	0.1533	0.2812	0.1312	0.0415	0.1799
shop_count					0.2261	0.2362	-0.0043	-0.0678	-0.0299	-0.1130
proc_count						<u>0.6993</u>	<u>0.4751</u>	0.2440	0.2053	<u>0.1923</u>
soci_count							<u>0.7595</u>	0.2239	0.1601	0.1177
sess_count								0.3821	0.2569	<u>0.3037</u>
xhour									0.8631	<u>0.2420</u>
yhour										0.1726

The ten features for respondent random3 that are most correlated with the target feature are displayed in table 4.1.3. The features that are most correlated with the target feature are noti_count, util_count, life_time and ente_time, with a correlation of respectively -0.1414, -0.1339, -0.1276 and -0.1238. There are also other features that are strongly related, these will be discussed in section 5.

	game_ time	life_ time	proc_ time	life_ count	prod_ count	util_ count	noti_ count	min_ len_ sess	if_ weekend	target
ente_time	-0.0738	-0.0305	<u>0.3659</u>	-0.0144	<u>0.2872</u>	0.2196	<u>0.3042</u>	0.2848	-0.0277	<u>-0.1238</u>
game_time		0.0431	0.2641	0.0704	0.1298	0.0920	0.0415	-0.0781	-0.1188	-0.1048
life_time			0.0405	<u>0.8917</u>	-0.0633	0.2917	0.2813	0.1227	0.1070	<u>-0.1276</u>
proc_time				0.0651	<u>0.5727</u>	<u>0.5074</u>	0.2268	0.0030	-0.1685	-0.0900
life_count					-0.0574	<u>0.3037</u>	0.2496	-0.0178	0.1737	-0.1060
prod_count						<u>0.6022</u>	0.1415	0.0957	-0.0968	-0.0973
util_count							<u>0.3036</u>	0.0082	-0.1760	<u>-0.1339</u>
noti_count								0.0257	-0.0623	<u>-0.1414</u>
min_len_sess									0.0871	-0.0967
if_weekend										-0.1090

The correlation for the features of the top30 respondents are displayed in table. There are no strong relations between the features and the target. The feature with the strongest relation with the target feature is sess_count with a correlation of 0.0842, which is very weak. Furthermore the total_time_phone and uni_app_count are strongly related with a portion of the other features.

Table 4.1.4 Top30 - 10 most influential features

	soci_ time	util_ time	game_ count	proc_ count	soci_ count	total_ time_ phone	sess_ count	min_ len_ sess	uni_ app_ count	target
game_time	0.0283	0.0009	<u>0.7569</u>	-0.0315	0.0242	0.1898	0.0706	-0.0218	0.1478	0.0517
soci_time		0.0668	0.0509	0.2306	<u>0.6032</u>	<u>0.7317</u>	0.1509	0.2515	<u>0.3800</u>	0.0399
util_time			-0.0088	0.2281	0.0206	<u>0.3594</u>	0.0609	0.0272	0.2229	-0.0433
game_count				-0.0217	0.0718	0.1581	0.1340	-0.0214	0.1585	0.0376
proc_count					0.3549	0.3247	0.2457	-0.0099	<u>0.5390</u>	<u>0.0675</u>
soci_count						<u>0.4516</u>	<u>0.5921</u>	0.0654	<u>0.5211</u>	<u>0.0643</u>
total_time_phone							0.1437	<u>0.3220</u>	<u>0.5138</u>	0.0423
sess_count								-0.1718	<u>0.5023</u>	<u>0.0842</u>
min_len_sess									0.0096	-0.0531
uni_app_count										0.1015

Uni_app_count, sess_count and soci_count are the features that have strongest correlation with the target feature. These relations are all very weak. The uni_app_counts shows a moderate relation with proc_count and soci_count, which is not strange since these applications are among the most used features. Furthermore the soci_time is strongly related with the total_time_phone, this is not remarkable as it is the most used application category.

Table 4.1.5 Topall - 10 most influential features

	soci_ time	life_ count	proc_ count	soci_ count	total_ time_ phone	sess_ count	min_ len_ sess	uni_ app_ count	xhour	target
life_time	0.0038	<u>0.5048</u>	0.0967	0.0312	0.2245	0.0940	0.0280	0.1698	0.0298	0.0590
soci_time		0.0268	0.2612	<u>0.6069</u>	<u>0.7003</u>	0.2243	0.2064	0.3914	0.0504	0.0419
life_count			0.1151	0.1084	0.1499	0.3338	0.0061	0.2707	0.0553	0.0440
proc_count				0.4700	0.3255	0.2938	-0.0001	<u>0.4531</u>	0.0310	0.0494
soci_count					0.4399	<u>0.6600</u>	0.0307	<u>0.5442</u>	0.1047	<u>0.0725</u>
total_time_phone						0.1876	0.2889	<u>0.5215</u>	0.0507	0.0412
sess_count							-0.1672	<u>0.5574</u>	0.1357	<u>0.0849</u>
min_len_sess								-0.0040	-0.0432	-0.0442
uni_app_count									0.0528	<u>0.0994</u>
xhour										0.0351

All the different models have their own specific features that have the most influence on the target feature. Sess_count belongs to the strongest related features for respondent random2 and this feature is also one of the stronger related features for the top30 and topAll models. This will be discussed in section 5.

4.1 Classification models

In this section the results of the second part of the research will be presented. Tests have been run to answer research question two: What model predicts best the perceived stress levels of smartphone users?

4.2.1 Random Forest

The random forest is one of the three classification algorithms that has been used in this study. The parameters of the random forest have been tuned with the use of the GridSearchCV instance. The parameters haven been tested with the values that are mentioned in section 3.8. The optimal parameters are different for each model and are displayed in table.

	N_ estimators	Max_ depth	Min_ samples_split	Min_ samples_leaf
Random1	15	25	4	1
Random2	20	50	5	3
Random3	8	15	4	2
Top30	30	10	2	1
TopAll	30	10	2	1

Four out of the five models achieved the highest accuracy with the Random Forest algorithm. Model Random1 has the best performance on the test set with the random forest with an accuracy of 94.1%, the baseline for this model was 81.7%. Model 1 scored high above the baseline. Also for model Random2 the random forest achieved the highest accuracy of 63.2% against a baseline of 63.8%. So it scored just below the baseline. The Top30 model scored an accuracy of 68.6% against a baseline of 69.6%, this model scored below the baseline. This is also the case for the TopAll model which achieved an accuracy of 67.1% against the baseline of 71.8%. Only Random1 model achieved an accuracy that was above the baseline. The table 4.2.2 shows the predictions of the random forest algorithm.

	Random1		Random2		Random3		Top30		Topall	
Actual/ Predicted	Class	Class	Class	Class	Class	Class	Class	Class	Class	Class
<u>Class 0</u>	<u>2</u>	1	<u>9</u>	3	<u>8</u>	4	<u>283</u>	58	<u>450</u>	96
<u>Class 1</u>	0	<u>14</u>	4	<u>3</u>	3	<u>1</u>	96	<u>53</u>	154	<u>61</u>

K-Nearest Neighbors (k-NN)

The second classification algorithm that was applied in this study, was the k-Nearest Neighbors. The data was normalized before analysis, since the k-NN algorithm requires normalized data. Since the K-NN algorithm is non-parametric, it is not difficult to find the optimal value of k. The different models are plotted and the optimal values of k can be noted from these plots. These optimal parameters can be found in table 4.2.3

	k
Random1	15
Random2	6
Random3	2
Top30	2
TopAll	2

For model random3 the k-NN algorithm scored the highest performance. This model had an accuracy of 75% which is just below the baseline of 76.3%. Model random1 scored 76.5% accuracy against the baseline of 81.7%. Model random2 scored only 52.5% accuracy against a baseline of 63.8%. Models Top30 and TopAll scored an accuracy of 60.6% and 62.6% against a baseline of 69.6% and 71.8% respectively. Table 4.2.4 shows the predictions of the k-Nearest Neighbours algorithm.

	Random1		Random2		Random3		Top30		Topall	
Actual/ Predicted	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1
Class 0	<u>2</u>	1	<u>7</u>	5	<u>11</u>	1	<u>260</u>	81	<u>415</u>	131
Class 1	3	<u>11</u>	4	<u>3</u>	3	<u>1</u>	112	<u>37</u>	157	<u>58</u>

Support Vector Machine Classifier

The last algorithm is the svm classifier. The parameters of this model are tuned with the use of the

	C	Epsilon
Random1	100	0.1
Random2	1	1
Random3	0.1	0.1
Top30	1	1
TopAll	1	0.1

GridSearchCV. The optimal parameters can be found in table 4.2.5. None of the model achieved a nice accuracy score with this algorithm. Table 4.2.6 provides the predictions of the SVM algorithm. The accuracy scores for the SVM algorithm on the train and test sets can be found in table 4.2.8 on the following page.

	Random1		Random2		Random3		Top30		Topall	
Actual/ Predicted	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1
Class 0	<u>2</u>	1	<u>6</u>	6	<u>5</u>	7	<u>203</u>	138	<u>316</u>	230
Class 1	1	<u>13</u>	3	<u>4</u>	2	<u>2</u>	67	<u>82</u>	91	<u>124</u>

(Balanced) Accuracy Test

Since the data was imbalanced the accuracy score was not very useful, the score did not give a representation of the performance. Therefore the balanced accuracy is used as an extra metric. As can be seen from table 4.2.7 show the balanced accuracy scores deteriorated scores. Only model1 achieves an balanced accuracy that is above the baseline on all of the algorithms. The rest of the models only achieves an accuracy below the baseline.

	RandomFores t	SVM	kNN	Baseline
Random1	96.7%	79.8%	65.8%	81.7%
Random2	59.6%	53.3%	50.6%	63.8%
Random3	46.4%	46.8%	64.3%	76.3%
Top30	61.2%	56.2%	50.6%	69.6%
TopAll	56.7%	56.3%	51.6%	71.8%

Table 4.2.8 Accuracy scores all models

Dataset	Baseline	Random Forest		Support Vector Machine		k-Nearest Neighbours	
		Train	Test	Train	Test	Train	Test
Random1	81.7%	100%	<u>94.1%</u>	98.1%	88.2%	80.2%	76.5%
Random2	63.8%	92.0%	<u>63.2%</u>	88.5%	52.6%	72.9%	52.6%
Random3	76.3%	96.9%	56.3%	67.3%	43.75%	98.0%	<u>75.0%</u>
Top30	69.6%	93.4%	<u>68.6%</u>	63.7%	58.2%	96.3%	60.6%
TopAll	71.8%	88.1%	<u>67.1%</u>	59.6%	57.8%	96.7%	62.2%

5. Discussion, Limitations & Conclusion

In this last section the discussion, limitations and conclusion can be found. Section 5.1 provides the discussion of the results. Section 5.2 discusses the limitations of the study, and Section 5.3 is the conclusion.

5.1 Discussion

The goal of this study was to analyse whether perceived stress levels could be predicted by phone application usage. This goal was split into two research questions, one that focused on the correlation of the features of the different models. The other part focussed on which model the best predictor was of the perceived stress levels of smartphone users.

The first part of the study focussed on the different features for all the different models. The aim was to explore for each model, if the stress level was strongly correlated with other features. Models random1, random2, and random3 all showed features that are moderately or strongly related with the stress features. These were smaller samples, which is also of influence but the different model also showed a lot of dissimilarities. For the three random models the most influential features were not the same. None of these three models had one of the others most influential features in their top three most influential features. This suggest that for every model other features are important. This is in line with results of another study (Ferdous, Osmani, & Mayora, 2015) (Bauer & Lukowicz, 2012) (Muaremi, Arnrich, & Tröster, 2013) that suggests that user-specific phone usage data is a better predictor for stress than generic data.

The second part of the study focussed on the predicting the stress with different algorithms. Only the model random1 achieved a balanced accuracy above the baseline on all the algorithms. The rest of the models scored below the baseline score. This suggests that the amount of data that is collected in the two hours before a survey might be too little. The performance of model Random1 on all the algorithms scored above the baseline. Model random1 seemed to be a lucky shot and is the exception on the rest. On average the algorithms scored badly. This might have to do with the minimum amount of data that is collected within those two hours upon a survey.

5.2 Limitations and future research

Two hour time frame

The decision to focus only one the two hours upon a survey might have been a little optimistic. The amount of data that is collected in this time frame is too minimal, as a result the observations have a high similarity. An idea for future research is the working with the moving time window technique. Which can be used with and without overlap.

Wasteful approach

The received survey dataset contained 16.016 entries, eventually 3.802 surveys are used for analysis. A large part of the dataset did contain empty surveys but this were only 5.372 surveys. Future researchers could make use of data imputation techniques to replace the NaNs.

The self-reporting bias

People that self-report their mood might be biased. These biases are categorised social desirability and recall bias. One might fill in on a survey that he thinks that he is stressed, but he only gives this stress level because he just almost had a car accident. With a survey this cannot be checked. For future research the application usage might be linked with a heartbeat sensor or other wearables that have physiological sensors.

5.3 Conclusion

The goal of this research was to analyse whether perceived stress levels could be predicted by phone application usage. To address this task there were two research questions developed:
Q1. Is there a distinction between the factors that influence the perceived stress levels for the user-specific models and the generic models?

Q2. What model predicts best the perceived stress levels of smartphone users?

The first question was answered through the use of Pearson correlaton matrices. Five different models were developed, 3 with user-specific data and two generic models of different amounts of data. All the user-specific models showed features that were moderately and strongly related with the stress level feature. But non of these features matched with the other models. This suggests that the impact that a phone has on the stress levels is different for each respondent. Which is in line with other studies. (Ferdous, Osmani, & Mayora, 2015) (Bauer & Lukowicz, 2012) (Muaremi, Arnrich, & Tröster, 2013)

Three different classification algorithms have been developed for the second research question: Random Forest, k-Nearest Neighbours and Support Vector Machine. The results showed that only one user-specific model scored above the baseline. The other models all scored below. The overall performance for the algorithms was poor.

The formulated problem statement for this research was:

To what extent can perceived stress levels be predicted by phone application usage?

The results showed that prediction of stress on basis of phone usage is not possible in the way that is has been addressed in this study. The models all scored below the baseline accuracy. This research has shown that the amount of data that is collected in the time frames before the survey is very little. In future research wider timeframes are suggested, that might even have overlap with other timeframes.

Bibliography

- Ahn, H., Wijaya, M. E., & Esmero, B. C. (2014). A systemic smartphone usage pattern analysis: focusing on smartphone addiction issue. *Int J Multimedia Ubiquitous Engineering*, 9-14.
- Altman, N. (1992). An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*, 46(3), 175-185.
- American Psychological Association. (2015). *Stress in America: Paying With Our Health*. Retrieved from APA: <https://www.apa.org/news/press/releases/stress/2014/stress-report.pdf>
- Andrews, S., Ellis, D., Shaw, H., & Piwek, L. (2015). Beyond Self-Report: Tools to Compare Estimated and Real-World Smartphone Use. *Plos One*.
- Annie, A. (2017). *Retrospective: A monumental year for the app economy*. App Annie.
- Bauer, G., & Lukowicz, P. (2012). Can Smartphone Detect Stress-Related Changes in Behaviour of Individuals. *IEEE International Conference on Pervasive Computing and Communications* (pp. 423 - 426). Lugano: IEEE.
- Beyan, C., & Fisher, R. (2015). Classifying imbalanced data sets using similarity based. *Pattern Recognition*, 48, 1653-1672.
- Blix, E., Perski, A., Berglund, H., & Savic, I. (2006). Long-term occupational stress is associated with regional reductions in brain tissue volumes. *PLoS One*, 8.
- Bogomolov, A., Lepri, B., Ferron, M., Pianesi, F., & Pentland, A. S. (2014). Pervasive stress recognition for sustainable living. *2014 IEEE International Conference on Pervasive Computing and Communication Workshops (PERCOM WORKSHOPS)* (pp. 345-350). IEEE.
- Bomhold, C. R. (2013, 3 29). Educational use of smart phone. *Program: electronic library and information systems*, pp. 424-436.
- Brefeld, U. (2019). *Machine Learning and Data Mining for Sports Analytics*. Springer.
- Breiman, L. (2001). *Random Forests*. Berkeley: University of California.
- Cawley, G. C., & Talbot, N. L. (2010). On over-fitting in model selection and subsequent selection bias in performance evaluation. *Journal of Machine Learning Research*, 11(Jul), 2079-2107.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: synthetic minority. *Journal of artificial intelligence research*, 321-357.
- Clark, A. (2018, 3 26). The Mind-Expanding Ideas of Andy Clark - The tools we use to help us think—from language to smartphones—may be part of thought itself. (L. MacFarquhar, Interviewer)
- Clarke, A., & Chalmers, D. (1998). The Extended Mind. *Analysis*, 7 - 19.

- Deloitte. (2019). *Global Mobile Consumer Survey 2019*. Deloitte.
- Epel, E. S., Crosswell, A. D., Mayer, S. E., Prather, A. A., Slavich, G. M., Puterman, E., & Mendes, W. B. (2018). More than a feeling: A unified view of stress measurement for population science. *Frontiers in neuroendocrinology*, *49*, 146-169.
- Evans, J. D. (1996). *Straightforward statistics for the behavioral sciences*. . Pacific Grove: Brooks/Cole Publishing.
- Ferdous, R., Osmani, V., & Mayora, O. (2015). Smartphone app usage a a predictor of perceived stress levels at workplace. *9th International Conference on Pervasive Computing Technologies for Healthcare (PervasiveHealth)*. Trento: CREATE-NET.
- Franke, S. (2003). Studying and working: The busy lives of students with paid employment. *Statistics Canada*.
- Fukazawa, Y., Ito, T., Okimura, T., Yamashita, Y., Maeda, T., & Ota, J. (2019). Predicting anxiety state using smartphone-based passive sensing. *Journal of biomedical informatics*, *93*.
- Fukunaga, K. (2013). *Introduction to statistical pattern recognition*. Elsevier.
- Godbout, J. P., & Glaser, R. (2006). Stress-induced immune dysregulation: implications for wound healing, infectious disease and cancer. . *Journal of Neuroimmune Pharmacology*, 421-427.
- Hamel, L. H. (2011). *Knowledge discovery with support vector machines (Vol. 3)*. John Wiley & Sons.
- Harri, A., & Brorsen, B. (2009). The Overlapping Data Problem. *Quantitative and Qualitative Analysis in Social Sciences*, *3*, pp. 78-115.
- Hendrickson, A., Abeele, M. V., & Aalbers, G. ((under review)). *Phone use data*.
- Hooftman, W., Mars, G., Janssen, B., de Vroome, E., Janssen, B., Pleijers, A., . . . van den Bossche, S. (2019). *Nationale Enquête Arbeidsomstandigheden 2018*. Leiden: TNO | CBS.
- Howlett, R. J., & Jain, L. C. (2001). Radial basis function networks 1: recent developments in theory and applications (Vol. 66). *Springer Science & Business Media*.
- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, *9(3)*, 90–95.
- James, G., Witten, D., Hastie, T., & R., T. (2017). *An Introduction to Statistical Learning*. New York: Springer.
- Jones, E., Oliphant, T., Peterson, P., & others, &. (2001). *SciPy: Open source scientific tools for Python*. Retrieved from <http://www.scipy.org/>
- Kang, J. M., Seo, S. S., & & Hong, J. W. (2011). Usage pattern analysis of smartphones . *13th Asia-Pacific Network Operations and Management Symposium* (pp. 1-8). IEEE.

- Kirschbaum, C. P. (1993). The 'Trier Social Stress Test'—a tool for investigating psychobiological stress responses in a laboratory setting. *Neuropsychobiology*, 76-81.
- Köhl, M., Magnussen, S. S., & Marchetti, M. (2006). *Sampling Methods, Remote Sensing and GIS Multiresource Forest Inventory*. Springer Science & Business Media.
- Kompier, M., & Cooper, C. (2003). *Preventing Stress, Improving Productivity - European case studies in the workplace*. London: Routledge.
- Kushlev, K., Proulx, J., & Dunn, E. (2016). "Silence Your Phones": Smartphone Notifications Increase. *CHI'16 - 2016 CHI Conference on Human Factors in Computing Systems* (pp. 1011 - 1020). San Jose: ACM.
- Larose, D. T., & Larose, C. D. (2014). *Discovering Knowledge in Data*. New Jersey: Wiley.
- Lebepe, F., Niezen, G., Hancke, G. P., & Ramotsoela, T. D. (2016). Wearable stress monitoring system using multiple sensors. *IEEE 14th International Conference on Industrial Informatics (INDIN)* (pp. 895-898). IEEE.
- Maimon, O., & Rokach, L. (2005). *Data mining and knowledge discovery handbook*. Springer.
- Matthews, K. A., Woodall, K. L., & Allen, M. T. (1993). Cardiovascular reactivity to stress predicts future blood pressure status. *Hypertension*, 22(4), 479-485.
- McKinney, W., Wes, & others, &. (2010). Data structures for statistical computing in python. *9th Python in Science Conference*, (pp. 51–56).
- Muaremi, A., Arnrich, B., & Tröster, G. (2013). Towards measuring stress with smartphones and wearable devices during workday and sleep. . *BioNanoScience*, 3, p. 172173.
- Oliphant, T. E. (2006). *A guide to NumPy (Vol.1)*. Trelgol Publishing USA.
- Oshiro, T. M., Perez, P. S., & Baranauskas, J. A. (2012). How many trees in a random forest? *International workshop on machine learning and data mining in pattern recognition* (pp. 154-168). Berlin: Springer.
- Oulasvirta, A., Rattenbury, T., Ma, L., & Raita, E. (2012). Habits make smartphone use more pervasive. *Personal and Ubiquitous Computing*, 16: 105-114.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., & ...others. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 2825 –2830.
- Pielot, M., & Rello, L. (2017). Productive, Anxious, Lonely - 24 Hours Without Push Notifications. *19th International Conference on Human-Computer Interaction with Mobile Devices and Services* (p. MobileHCI '17). Austria, Vienna: ACM.
- Przybylski, A. K., & Weinstein, N. (2017). A large-scale test of the Goldilocks Hypothesis: Quantifying the relations between digital-screen use and the mental well-being of adolescents. *Psychological Science*, 204-215.

- Rijksoverheid. (2019, 7). *Mag ik appen, bellen en naar muziek luisteren op de fiets?* Retrieved from Rijksoverheid.nl: <https://www.rijksoverheid.nl/onderwerpen/fiets/vraag-en-antwoord/mag-ik-bellen-en-naar-muziek-luisteren-op-de-fiets>
- Rosengren, A., Hawken, S., Ôunpuu, S., Sliwa, K., Zubaid, M., Almahmeed, W. A., & investigators, .. &. (2004). Association of psychosocial risk factors with risk of acute myocardial infarction in 11 119 cases and 13 648 controls from 52 countries (the INTERHEART study): case-control study. *The Lancet*, *364*, 953-962.
- Sarker, I. H., Kayes, A. S., & Watters, P. (2019). Effectiveness analysis of machine learning classification models for predicting personalized context-aware smartphone usage. *Journal of Big Data* *6*(1).
- Shu, X. (2020). *Knowledge Discovery in the Social Sciences: A Data Mining Approach*. University of California Press.
- Simyo, DirectResearch. (2019, 7). *Word jij slimmer of juist socialer van de smartphone?* Retrieved from Simyo: <https://www.simyo.nl/onderzoek/smartphone-gebruik/>
- Tan, P., Steinbach, M., & Kumar, V. (2014). *Introduction to Data Mining*. Essex: Pearson.
- Thoméé, S. (2018). Mobile Phone Use and Mental Health. A Review of the Research That Take a Psychological Perspective on Exposure. *International Journal of Environmental Research and Public Health*, *15*.
- Waskom, M., Botvinnik, O., O'Kane, D., Hobson, P., Ostblom, J., & Lukauskas, S. (2018, 07). *mwaskom/seaborn: v0.9.0*. Retrieved from <https://doi.org/10.5281/zenodo.1313201>
- Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann.
- Zenonos, A., Khan, A., Kalogridis, G., Vatsikas, S., Lewis, T., & Sooriyabandara, M. (2016). HealthyOffice: Mood recognition at work using smartphones and wearable sensors. *2016 IEEE International Conference on Pervasive Computing and Communication Workshops (PerCom Workshops)* (pp. 1-6). IEEE.

Appendix B

3.2.1 Duplicate data

The phone dataset consists of 586.792 entries, in this data there are 19.280 erroneous duplicate entries. The mood survey dataset contained 14 duplicate entries. These duplicate entries are removed from the data.

3.2.2 Erroneous User ID's

The mood survey dataset consists of 16.002 entries from 149 different users. There are 124 unique user IDs in the phone usage dataset and 149 unique user IDs in the mood survey dataset. This indicates that there are 25 erroneous user IDs in the mood survey dataset. These 25 user IDs are responsible for another 1.949 surveys. These surveys are excluded from the data. This reduced the sample size from $n_t = 16.002$ to $n_t = 14.053$.

3.2.3 Expired, double, blocked and cancelled surveys

The survey data also consists of non-filled and erroneous filled-in surveys. The moment that the respondents received the surveys, they were given two hours to complete it. After this moment the surveys would expire. The dataset consists of 5.372 returned expired surveys that are empty and 33 expired surveys that are partly filled-in. The entries that are partly filled-in miss multiple data and are not workable for analysis. Furthermore the respondents were also given the possibility to block or cancel a survey. There are 128 surveys in the dataset that are blocked or cancelled. All these expired, cancelled and blocked surveys are removed from the dataset. This reduced the sample size from $n_t = 14.053$ to $n_t = 8.520$.

3.2.4 Failing software

Six of the respondents in the study have not returned a single survey while their phone was being logged. It is unclear what the cause of this is, but it seems that it occurred due to failing software. These users have sent 241 surveys, these surveys are omitted from the data.

Due to incorrect settings some of the participants also already send surveys before the phone was logged, but also kept receiving and replying surveys after the study was over. A total amount of 1.026 surveys was due to this useless. This brings the total amount of surveys from $n_t = 8.520$ to $n_t = 7.253$.

Another 29 respondents have returned surveys with the Ethica application at times that their phone was logged but the MobileDNA application did not register telephone use. It suggests that at least either of the two applications has malfunctioned. It was impossible to retrieve

whether the applications only malfunctioned at these moments or during the whole study period, therefore these respondents are excluded from the data. These respondents are responsible for 930 surveys, which brings the total amount of surveys from 7.253 to 6.323.

3.2.5 Outliers

Two applications were excluded from the data and placed in the ‘wrong’ category. One of these applications is ‘com.madein.coinpotapp’, which is an application that mines bitcoins. This application runs on the background and is not actively used by the phone user but seen as actively used. Only one respondent had this applications installed. The other application is ‘ru.woxTGdZL.IfEhxxjCz’ which is an unknown application that is used by only one user for a limited amount of times(n=5). The two applications that are placed in the wrong category are will not be included in the data for future analysis.