

Technische Universität Berlin

Institut für Softwaretechnik und Theoretische Informatik
Quality and Usability Lab

Fakultät IV
Franklinstrasse 28-29
10587 Berlin



Master Thesis

Predicting tap locations on touch screens in the field using accelerometer and gyroscope sensor readings

Emanuel Schmitt

Matriculation Number: 333772

Examined by:
Prof. Dr.-Ing. Sebastian Möller
Prof. Dr.-Ing. Axel Kuüper

Supervised by:
Dr.-Ing. Jan-Niklas Antons

Thanks to all dem brothas

Hereby I declare that I wrote this thesis myself with the help of no more than the mentioned literature and auxiliary means.

Berlin, 01.01.2050

.....
(*Signature [your name]*)

Abstract

This template is intended to give an introduction of how to write diploma and master thesis at the chair 'Architektur der Vermittlungsknoten' of the Technische Universität Berlin. Please don't use the term 'Technical University' in your thesis because this is a proper name.

On the one hand this PDF should give a guidance to people who will soon start to write their thesis. The overall structure is explained by examples. On the other hand this text is provided as a collection of LaTeX files that can be used as a template for a new thesis. Feel free to edit the design.

It is highly recommended to write your thesis with LaTeX. I prefer to use MikTeX in combination with TeXnicCenter (both freeware) but you can use any other LaTeX software as well. For managing the references I use the open-source tool jabref. For diagrams and graphs I tend to use MS Visio with PDF plugin. Images look much better when saved as vector images. For logos and 'external' images use JPG or PNG. In your thesis you should try to explain as much as possible with the help of images.

The abstract is the most important part of your thesis. Take your time to write it as good as possible. Abstract should have no more than one page. It is normal to rewrite the abstract again and again, so probably you won't write the final abstract before the last week of due-date. Before submitting your thesis you should give at least the abstract, the introduction and the conclusion to a native english speaker. It is likely that almost no one will read your thesis as a whole but most people will read the abstract, the introduction and the conclusion.

Start with some introductory lines, followed by some words why your topic is relevant and why your solution is needed concluding with 'what I have done'. Don't use too many buzzwords. The abstract may also be read by people who are not familiar with your topic.

Zusammenfassung

Contents

List of Figures	xiii
List of Tables	xv
1 Introduction	1
1.1 Motivation	1
1.2 Objective	1
1.3 Outline	1
2 Related Work	3
2.1 Eavesdropping on Emanations	3
2.1.1 Acoustic Emanations	3
2.1.2 Optical Emanation	4
2.1.3 Electro-magnetic Emanation	5
2.1.4 Motion Emanation	5
2.2 Eavesdropping on touch screen user interactions	6
3 Machine Learning Fundamentals	11
3.1 Overview and Definition	11
3.2 Categorization of Methods	12
3.3 Support Vector Machines	14
4 Implementation	17
4.1 System Architecture	17
4.2 Mobile Application	18
4.2.1 Interface	18
4.2.2 Smartphone Sensors	19
4.3 Backend application	21
4.3.1 General requirements	21
4.3.3 Data model	21
4.3.4 User management & Security	21
4.3.5 Ensuring consistency	21
5 Inferring Tap Locations	23
5.1 Data Preprocessing	23
5.2 Convolution Neural Network	23
5.3 Recurrent Neural Network	23

6	Method	25
6.1	Overview	25
6.2	Hypothesis	25
6.3	Experimental Setup	25
6.3.1	Subjects	25
6.3.2	Devices	25
6.3.3	Experiment Settings	25
6.4	Analysis	25
7	Results	27
7.1	Hypothesis	27
7.2	Discussion	27
8	Conclusion	29
8.1	Further Outlook	29
	List of Acronyms	31
	Bibliography	33
	Annex	37

List of Figures

4.1	TapSensing architecture overview.	17
4.2	Grid sizes of Grid.	18
4.3	The figure displays question views with icons as answer possibilities. . . .	19
4.4	Apple iPhone with the corresponding axes of the accelerometer.	20
4.5	Swift code snippet displaying how to access sensor values with Core Motion.	21

List of Tables

1 Introduction

This chapter should have about 4-8 pages and at least one image, describing your topic and your concept. Usually the introduction chapter is separated into subsections like 'motivation', 'objective', 'scope' and 'outline'.

1.1 Motivation

Start describing the situation as it is today or as it has been during the last years. 'Over the last few years there has been a tendency... In recent years...'. The introduction should make people aware of the problem that you are trying to solve with your concept, respectively implementation. Don't start with 'In my thesis I will implement X'.

1.2 Objective

1.3 Outline

The 'structure' or 'outline' section gives a brief introduction into the main chapters of your work. Write 2-5 lines about each chapter. Usually diploma thesis are separated into 6-8 main chapters.

This example thesis is separated into 7 chapters.

Chapter 2 is usually termed 'Related Work', 'State of the Art' or 'Fundamentals'. Here you will describe relevant technologies and standards related to your topic. What did other scientists propose regarding your topic? This chapter makes about 20-30 percent of the complete thesis.

Chapter 4 analyzes the requirements for your component. This chapter will have 5-10 pages.

Chapter 5 is usually termed 'Concept', 'Design' or 'Model'. Here you describe your approach, give a high-level description to the architectural structure and to the single components that your solution consists of. Use structured images and UML diagrams for explanation. This chapter will have a volume of 20-30 percent of your thesis.

Chapter 6 describes the implementation part of your work. Don't explain every code detail but emphasize important aspects of your implementation. This chapter will have

a volume of 15-20 percent of your thesis.

Chapter 7 is usually termed 'Evaluation' or 'Validation'. How did you test it? In which environment? How does it scale? Measurements, tests, screenshots. This chapter will have a volume of 10-15 percent of your thesis.

Chapter 8 summarizes the thesis, describes the problems that occurred and gives an outlook about future work. Should have about 4-6 pages.

2 Related Work

The utilization of motion sensors sidechannels in order to reconstruct user interactions can be seen as a form of eavesdropping. Hence, in order to shed light into this category of security breaches, this chapter will cover past research based on the usage of sidechannels emanations which can reveal confidential information. Therefore, The first section will cover different techniques that utilize acoustic, electromagnetic, optical and motion emanation for eavesdropping purposes while the second section shows related work aiming at inferring tap interactions on touchscreens.

2.1 Eavesdropping on Emanations

Eavesdropping is defined as the the practice of secretly listening to private conversations of others without their consent [6]. As this definition originally refers to conversations between humans, eavesdropping can also be seen in terms of human-computer-interaction. In this context, the computer is seen as one conversational partner whereas the user interacting with the device is seen as the other. As the channel of communication in a human conversation is the acoustic channel, the channel in which human-computer interaction takes place has various forms. Typically, a human may enter information on a peripheral device while the computer gives feedback through an image representation. However, speech interfaces and other forms of interaction are also possible. In order to spy on the human-computer conversation, a third party with malicious intent must apply different techniques in order to spy on these interactions.

In this context, a frequently discussed practice throughout academia involves the use of leaking emanations for eavesdropping purposes. These emanations can be monitored in order to carefully reconstruct the contained information. As these emanations occur in various formats, a categorization has been done based on the sensory channel they transpire. Therefore, the following sections will cover eavesdropping techniques based on acoustic, optical, electromagnetic and motion emanations.

2.1.1 Acoustic Emanations

One way of spying on electrical devices is by utilizing the acoustic channel [4, 1, 40]. Many electronic devices deploy tiny mechanics that generate sounds as a byproduct during interactions or during operation. These distinct sounds can differ in their characteristics making them adequate to identify the original information currently being processed by the machine. In these scenarios an eavesdropper targets a microphone

in near proximity of the target device to capture the audio signals the device is exposing. A learning algorithm is then applied to the audio signals to reconstruct information.

In 2010, Backes et al. examined the problem of acoustic emanations of dot matrix printers, which where, at that time, still commonly used in banks and medical offices. By using a simple consumer-grade microphone, the researchers were able to recover whole sentences the printer was printing based. Backes et al. processed the audio samples in order to extract frequency-domain features. These features worked as input for a hidden markov model, a technology commonly used in audio speech recognition. As a result, the recognition system was able to reconstruct individual characters based on the sound inputs. To demonstrate a potential attack, the researchers deployed the system in a medical office being able to obtain up to 72% of the sentences being printed on medical subscriptions [4].

Being inspired by the findings concerning the dot matrix printer, Asonov and Agrawal investigated acoustic emanations produced by hitting keystrokes on a desktop and a notebook keyboard. Following their hypothesis claiming that each keystroke has a macroscopic difference in its construction mechanics, as well as a distinct reverberation caused by the position in the board, individual keystrokes were recorded [1]. Researchers then extracted frequency domain features from the audio signals and passed them into a neural network. In an experiment performing 300 keystrokes, 79% of the characters could be correctly recognized [1]. As this technique required substantial training before recognition, other studies have reached similar accuracies using an unsupervised approach [40] on the one hand and by using acoustic dictionaries [5] on the other.

2.1.2 Optical Emanation

Besides acoustic emanations, optical emanations can also pose a valuable source of information for a potential eavesdropper. Most electronic devices, such as notebook, smartphones and tablet computers, provide graphical user interfaces through their own built-in screens. Even though these screens are meant to target the human eye, they can reflect off other surfaces. The reflections can be caught by high resolution camera sensors, which can then display the image revealing secret information.

One example of the use of optical emanations has been developed by Kuhn aiming to eavesdrop on cathode-ray-tube(CRT) monitors at distance. The researcher has shown that the information displayed on the monitor can be reconstructed from its distorted or even diffusely reflected light. In an experiment, Kuhn targeted a screen displaying an image against a wall while the reflections of the screen were captured using a photomultiplier. The experiment showed that enough high-frequency content remained in the emitted light for a computer to reconstruct the original image [16]. A similar approach that comprises reflections has been shown by Backes et al., however focusing on LCD displays. In this experiment, the researchers caught reflections in various objects that are commonly to be found in close proximity to a computer screen. Such objects included

eyeglasses, tea pots, spoons and even plastic bottles. This work was later extended to additionally capture screens based on the reflections on the human eye's cornea [3].

2.1.3 Electro-magnetic Emanation

As electric currents flow through computer components, they emit electromagnetic waves to their near surrounding. These electromagnetic radiations can be picked up as a side channel using sensitive equipment in order to retrieve data. Electromagnetic emanations have been a research topic that has been present for several decades while the first research conducted reaches far in time.

Back in 1943, a research group under the codename TEMPEST, a subdivision of the NSA¹, were able to infer information from the infamous Bell Telephone model 131-B2, a teletype terminal which was used for encrypting wartime communication [29]. Using an oscilloscope, researchers could capture leaking electromagnetic signals from the device and by carefully examining the peaks of the recorded signals, the plain message the device was currently processing could be reconstructed [29]. This technique was later advanced and used in the Vietnam war. Through similar electric emanation the US military could detect approaching Viet Cong trucks giving them an immense competitive advantage [26]. Today, TEMPEST is a security standard for electronic devices ensuring that certified devices do not accidentally emanate confidential information [27].

A second prominent finding of eavesdropping on electromagnetic emanations was done by the researcher van Eck, who discovered that cathode-ray-tube monitors could be spied upon from a distance [37]. By using general market equipment, such as antennas van Eck and standard receivers, signals emitted from the cable connecting the computer to the monitor could be received. Since these cables only transmit the video signal for visualization, the researcher could display the visual output of the target monitor revealing a full screen cast of the original image. This attack is referred to in literature as *Van Eck Phreaking* [11, 15].

Furthermore, research has shown that side-band electromagnetic emanations are present in keyboards [38], computer screens [37, 17], printers [32], computer interfaces, such as USB 2 [28] and the parallel port [36] and in Smart Cards[33].

2.1.4 Motion Emanation

In the past decade modern devices are increasingly equipped with highly responsive sensors, such as the gyroscope and accelerometer enabling the devices to sense rich interactions with their environment. As user interactions, such as typing the keyboard or tapping on touchscreens, require the user to apply a certain force while entering information, this motion can be captured by motion sensors in order to be used for a

¹National Security Agency

side-channel attack [24, 30, 7].

Marquardt et al. conducted an experiment where an Apple iPhone that captures accelerometer motion was placed next to a desktop keyboard. Subjects then had to enter sentences while the application was monitoring the user's motion. The researchers could decode the accelerometer signals by mapping the certain vibration caused the typing motion to their keystrokes. The decoded characters were then matched based on a dictionary containing a frequency distribution of commonly used words. As a result, words could be successfully obtained with an accuracy of up to 80% [22].

As motion emanations are highly relevant for the work in this thesis, the next section will be dedicated to further work regarding this topic.

2.2 Eavesdropping on touch screen user interactions

As we have seen in the previous section, user interactions with peripheral devices, such as the keyboard or PIN pads, can be obtained by either the acoustic channel, by electromagnetic leakage or by capturing the motion of the user. However, with the rise in soft keyboard usage, the same discussed methods that extract information from keyboard do not apply to tap interactions on a touchscreen surface. Since a touchscreen does not embody fine mechanics producing sounds nor does it have emanating cables, the research community has developed nouvelle methods to spy on user inputs based on the smartphone embodied motion sensors.

The general idea behind the three approaches that are going to be discuss in the following is that a tap, or to be more precise the magnitude of the force of a tap, on a specific touch screen location creates an identifiable pattern on the motion sensors that can be sufficient to infer the initial tap location. This is particularly interesting since motion sensors are not considered as being privacy-sensitive and therefore lack access restrictions by the operating systems of the devices.

Touchlogger

The first paper regarding this security threat was published by Cai and Chen. In their proof-of-concept study they created an Android² application which displays a 10-digit PIN-pad. During interactions with the PIN-pad, the accelerometer signals were monitored and used for later data analysis. Having observed that a tap movement affects the rotation angle of the screen, the researchers handcrafted features based on the path of the *pitch*³ and the *roll*⁴ angles of the accelerometer. These were intersected to find

²Android operating system for smartphones by Google Inc.

³The pitch-angle corresponds to the x-axis of the accelerometer.

⁴The roll angle corresponds to the y-axis of the accelerometer.

a dominating edge on where the tap had presumably taken place. By using a probability density function for a Gaussian distribution the researchers were able to achieve an average accuracy of 70% for interred PIN-pad digits. The training set size involved 449 pin strokes [7]. Even though *Touchlogger* was a promising first step, due to its low granularity of only 10 distinguishable large screen areas, it remained unclear if the attack can be carried over to a full software keyboard. Furthermore, since the inference was performed on only a single smartphone model, the question is left open whether other smartphones or tablet computers are similarly vulnerable.

ACCessory

In order to show the feasibility for a full software keyboard, Owusu et al. performed a second attempt to the problem by creating *ACCessory*. *ACCessory* is an Android application with functionalities similar to the previously mentioned *Touchlogger*. However, the application significantly differ in its tap area granularity providing two separate modes for tap inputs: *area mode* and *character mode*. *area mode* consists of tap areas arranged in a 60-cell grid, whereas a QWERTY keyboard within landscape orientation was displayed in *character mode*. Having extracted features mainly from the time-domain, a classification using the Random Forests algorithm reached an accuracy of 24.5% for the 60-cell grid. Here, the corresponding dataset consisted of 1300 keystrokes collected from 4 participants. As the *area mode* experiment focused on recognizing individual keystrokes, the *character mode* experiment focused on cracking passwords. By combining the keystrokes into a sequence and assuming recognition errors in individual characters, the researchers could create a ranked list of candidate passwords by running a maximum likelihood search for the most probable classification errors for an obtained password. Here, 6 out of 99 password could be inferred under 4.5 median trials given that one trials refers to traversing down one item of the candidate list. Furthermore, the majority of 59 out of 99 passwords could be inferred within 2^{15} median trials. As general result, even though the overall accuracy of the learning system scored low, the researchers could significantly reduce the search space for reconstructing a password indicating that accelerometer readings can in deed yield confidential information.

TapPrints

The most comprehensive study to date regarding the topic was conducted by Miluzzo et al. and differs from *ACCessory* and *Touchlogger* in many important ways. While both previous studies are both evaluated on the Android smartphones, *TapPrints* investigates the tap inference on both iOS and Android operating systems including tablets and smartphones alike. Another important point of differentiation is the used learning system. In order to raise the level of entropy, *TapPrints* combines readings from the accelerometer and gyroscope for a more sophisticated feature extraction. Here, time-domain and frequency-domain features, as well as the correlation and angles between individual sensor components are considered. For classification purposes, the researchers use an ensemble method combining decision tress, support vector machines, k-nearest

neighbors and multinomial logistic regression in a winner-takes-it-all⁵ voting fashion. The dataset collected in this experiments contains over 40.000 individual taps collected from 10 different users. In addition, the researchers also requested user to use different input modalities while typing including the usage of the index finger and thumb.

The *TapPrints* undertaking consists of two separate experiments: The first is a icon tapping experiment where icons are arranged in a 20 cell grid and the second one being a letter tapping experiment involving the standard software keyboard offered by the operating system. In the first experiment, an average accuracy of 78% was achieved for the iPhone whereas 67% of icon taps could be correctly inferred on the Android device. In the letter tapping experiment users were asked to enter pangram⁶ sentences on the OS soft-keyboard. Results showed that on both iPhone and Android an average of 43% of the letters could be correctly classified. Even though the average accuracy for individual letters were not particularly high, Miluzzo et al. could show that when pangrams were repeatedly entered, a majority vote could be applies to individual character recognitions allowing to recover the whole pangram in approx. 15 trials. To conclude, *TapPrints* could demonstrate that motion sensor can be used to obtain passwords on multiple platforms and formats and with different input modalities.

Comparison to similar studies

Since the data used in *TapPrints* and in the other related work was collected in a controlled environment [30, 7, 24], it is not possible to tell if the feasibility of tap inference will also apply to data collected in a field environment. As this has not been investigated, this will form the central research question in this thesis. Additional questions will be discussed in the hypothesis section.

To draw the border between this study and the previous mentioned ones, this study will differ or relate as follows:

- **User interface:** For data acquisition, the user interface will cover tap area grids, as we have seen in *ACCessory* and *Touchlogger*, with 4, 12 and 20 distinguishable classes.
- **Devices:** Unlike *Touchlogger* and *ACCessory*, the devices used for this study are on the iOS Platform. Apple iPhone 6, 6s and 7 will be considered due to their mutual screen size.
- **Motion sensors:** *TapPrints* has shown that the gyroscope yields more information than the accelerometer [24], therefore both sensors will be monitored. Furthermore, sensors will be read at a frequency of 100Hz, since this had been proven to deliver best results [24, 30].

⁵This implies that all classifiers classify separately and the classifier with the highest prediction score wins.

⁶A pangram is a sentence using every letter of a given alphabet at least once.

- **Dataset:** In comparison to related studies, a dataset containing data points collected in a laboratory environment and the field environment will be acquired from a total of 27 users. This data will contain the index finger and thumb as input modalities and standing and sitting as body postures while input.
- **Feature extraction:** As *TapPrints* shows a deliberate list of features [24], that will be partially adopted. All features will be covered except features concerning the frequency domain. A more detailed look will be covered in section Feature Extraction
- **Classification:** A feed-forward neural network, as well as a support vector machine with radial kernel will be used for classification. More will be covered in the classification section.

3 Machine Learning Fundamentals

As machine learning techniques will be used for the classification of individual tap locations on a smartphone touchscreen, the following chapter will give a brief overview of the fundamental concepts evolving around statistical learning. Furthermore, the methods that going to be applied to the acquired sensor data will be discussed.

3.1 Overview and Definition

Ever since computers were invented, there has been a desire to enable them to learn [35]. This desire has grown into the field of machine learning which seeks to answer questions on how to build systems that automatically improve with experience, and what the fundamental laws of learning processes are. Today, state-of-the-art ML covers a large set of methods and algorithms designed to accomplish tasks where conventional hard-coded routines have brought insufficient results. From speech recognition to email spam detection or recommendation systems, ML methods find broad usage in a variety of problem domains.

In order to understand what the principle of machine learning is, we will start with a definition by Samuel [35]:

Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.

In this definition, special emphasis is to be put on the last part of this definition. A computer is only then able to learn when he can perform a task without being explicitly instructed. Thus, in order to learn, the computer must somehow be able to instruct itself without the influence of an outer . As this definition lacks a more detailed view on what computer learning is, we will dive into a definition by Tom Mitchell [25]:

A learning system is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

The example that Mitchell notes, is one from the games of checkers [25]. In this case checkers is the task T that the computer is aiming to learn. In order for the computer

to learn, information on previously played matches is required. Since the computer does not know how to evaluate if a particular match was either good or bad, we set the performance to be defined based on how many matches were actually won. If a computer program can raise the amount of games won (*performance measure P*) with the help of the experience from previously played matches (*experience E*) then it can learn to play checkers (*task T*).

To break this down into a more practical perspective, the challenge lies in finding an appropriate model in order to learn from data, which is the most common format to represent past experience. By learning, the computer adjusts parameters on the model based on the data that we feed the system with. Once the model has been adjusted, it can perform tasks with new incoming experience.

3.2 Categorization of Methods

As machine learning algorithms and methods differ from their approach to learning and underlying concepts, it is common practice to separate these into the following categories [9, 23] : Supervised learning, unsupervised learning, reinforcement learning and evolutionary learning. In the following sections I will briefly outline these.

Supervised learning, which is also named learning from example, is presumably the most prominent category of ML algorithms. The algorithm is given a training set of examples $\{x_0, \dots, x_n\}$, which are also known as *features* and the correct target values $\{y_0, \dots, y_n\}$ mapped to each set of features, which is the answer that the algorithm should produce. The algorithm then generalizes based on the training set in order to respond with sensible outputs on all possible input values. Outputs, if they are discrete labels, correspond to a classification task whereas outputs on a continuous scale refer to a regression task (see [23]).

An example for supervised learning is the classification of malignant or benign tumors as seen in cancer diagnosis. Let's assume we have a dataset with different properties of a tumor, such as the size or the color of the cells. These properties form our features x . Each set of features is mapped to an output label y stating if the tumor is malignant or benign. The first step is to use the pairs (x, y) of the training set to teach the algorithm the correct mapping of the problem space. As x is linked to the output y in the training set, learning the conjunction of these two values is done under supervision since the output label y is given. Once learned, the algorithm is generalized to map unseen inputs to the correct output label.

Practical applications are for example digit and handwriting recognition [18], spam filtering [10] for e-mails or network anomaly detection [19].

Presumably the most widely known machine learning techniques belong to this category, such as Support Vector Machines (SVMs), Artificial Neural Networks,

Unsupervised learning is the task of learning structures from input values that are not explicitly labeled. In comparison to supervised learning, where correct output values are provided for each input, unsupervised algorithms learn to identify similarities in the input data and can therefore group these [9]. These grouping problems are referred to as *clustering*. The underlying idea here, is that humans learn by not explicitly being told what the right answer should be [23]. If a human sees different species of snakes, for instance, he or she is able to identify them all as snakes. Hence, the human is aware that there are differences in each specific type of snake without specifically knowing a correct label.

A prominent example where unsupervised learning is heavily used, is in recommender systems for online retail shops. Amazon.com, for instance, uses a technique called *collaborative filtering*, which measures similarity in customers based that they have previously bought [20]. Having identified similar customers utilizing the cosine similarity, the algorithm can then recommend items that similar users have bought. This technique is also used for music recommendations [31] or social network recommendations [14].

The field of unsupervised learning is closely related to density estimation in statistics, as with the density of inputs, we are able to group them. The K-means algorithm is the most prominent in this field [23].

Reinforcement learning falls in between supervised and unsupervised learning methods. Whereas supervised learning tries to bridge the gap between input and corresponding output values and unsupervised methods detect groupings in incoming data, reinforcement learning is based on learning with a *critic* [23]. The algorithm tries different solution strategies to a problem and is told whether or not the answer provided was correct. An important fact here, is that the algorithm is not told how to correct itself. This practice of "trying-out" is based on the concept of *trial-and-error learning* which is known as the *Law-of-effect* [23]. A good example is a child that tries to stand up and learn walking. The child tries out many different strategies for staying upright and receives feedback from the field based on how long it can stand without falling down again. The method that previously worked best is then repeated in order to find the optimal solution resulting in the child learning to walk [23].

In more mathematical terms, the reinforcement learning problem is formalized with an agent and his environment. The environment in which the agent is set provides a set of *states* on which the agent can perform *actions* to maximize a certain *reward*. By performing actions the state changes and a new reward is calculated. The reward then tells the agent if the action was a good choice. Goal of the algorithm is to maximize the reward [23].

Reinforcement learning is a practical computational tool for constructing autonomous systems that improve themselves with experience. These applications have ranged from robotics, to industrial manufacturing, to combinatorial search problems such as computer game playing [13]. Prominent methods of this category are Q-learning, Monte Carlo methods and Hidden Markov Models [23].

Evolutionary learning is inspired by strongly inspired by nature. As biological evolution improves the survival of a species, the strategy of adaptation to improve survival rates and the chance of offspring has inspired researchers to craft genetic algorithms (GA) [23].

Genetic algorithms are a family of adaptive search procedures which have derived their name from the fact that they are based on models of genetic change in a population of individuals. These models have their foundation in three basic ideas: (1) Each evolutionary state of a population can be evaluated on a *fitness* scale. This is done since biological evolution has a natural bias towards animals that are fitter than others. These animals tend to live longer, are more attractive and generate healthier and happier offspring, an idea which was originated in Charles Darwin's *The Origin of Species*; (2) Each population can be mated to generate offspring using a *mating operator*. (3) The third component are *genetic operator*, such as *crossover* and *mutation*, which determine how the offspring solution is composed of the genetic material of the parents [8].

Evolutionary learning is often considered when other methods fail to find a reasonable answer. Algorithms find applications in search and mathematical optimization, but also in arts and simulation [23].

In this section we have seen several different problems that we can solve with the help of algorithmic learning. For our use case, as we want to predict the locations on smartphone screens using sensory data. As this is a supervised learning problem, we will cover one supervised approach in more detail in the following section: Artificial neural networks.

3.3 Support Vector Machines

A SVM is a non-linear kernel based extension of the so-called maximum margin classifier. Originating from a binary classification problem, where $y \in \{1, -1\}$, the general idea of a maximum margin classifier is to find a separating hyperplane in the p -dimensional feature space that separates the training instances which leads to a maximum distance between the observations of the two classes (James et al. 2013). This distance is referred to as margin M measuring the smallest distance of a training observation towards the defined hyperplane. In order to find a suitable hyperplane only a subset of the training examples are sufficient.

Despite all that, linear separability is unlikely to be applicable to most real world problems. However, when a problem is linear separable, the constraint to the support vectors induce an undesirable sensitivity to individual observations, which can lead to high variance of the classifier. For this very reason, the maximum margin classifier can be extended to it's more robust successor, the support vector classifier. It is still a linear classifier however allowing for violations in fitting the margin.

Mathematically, the support vector classifier can be described via following optimization problem:

4 Implementation

For labeled data acquisition a system was required that is able to function in a laboratory environment, as well as receive the data coming from the field study. For this purpose TapSensing was created. TapSensing is an iOS application that collects touch events including their sensory information to then send them to a backend server application. In the following sections we will outline the different components of TapSensing.

4.1 System Architecture

The TapSensing application comprises two main components, one being the mobile application and the other being the server-side application. The mobile client is responsible for generating the sensor and tap information by providing a user interface for the user to tap. For the data to be stored in a centralized manner, the backend provides endpoints as a gateway to the database.

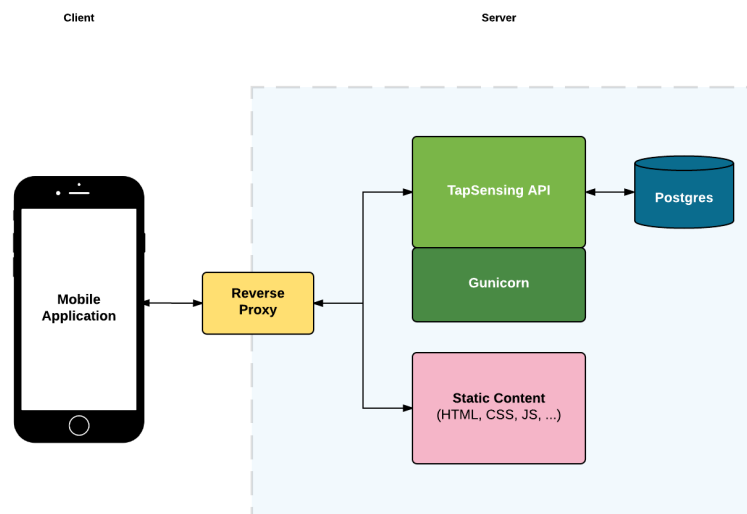


Figure 4.1: TapSensing architecture overview.

For the network requests containing the tap information which come from the mobile device to reach the backend, it must first pass through a reverse proxy. As reverse proxy we have chosen NGINX due to its easy configurability. In this case, NGINX forwards requests to the TapSensing application and serves static files.

The TapSensing backend is written upon the Python Django Framework¹ which is being executed upon the gunicorn application server. Django uses a so-called ORM to perform transactions with the Database, which in our case is a PostgreSQL database.

4.2 Mobile Application

4.2.1 Interface

Tap Input

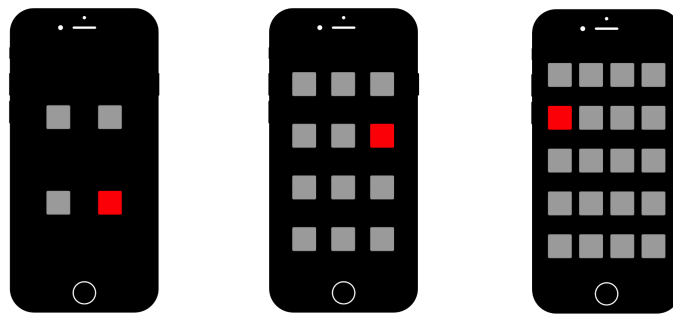


Figure 4.2: Grid sizes of Grid.

For the user to be able to enter the taps into the device the app provides a screen with buttons aligned in grid structures.

The positions of the buttons are calculated based on the configuration that is set. Here, the amount of buttons in height and width can be adjusted.

The source code of the of the question view is to be found in `GridViewController.swift`.

Questions

To obtain more information on the taps provided in the tap input view, the application provides views for the user to answer several questions concerning input modalities, body posture and mood. After the tap input screen, the questions are displayed providing multiple answer choices. In addition, the application provides a pictogram for each answer option.

Each view is generically set based on the questions and answer possibilities. Once the user taps on an icon the view transitions to the next question view.

The source code of the of the question view is to be found in `QuestionViewController.swift`.

¹<https://www.djangoproject.com/>



Figure 4.3: The figure displays question views with icons as answer possibilities.

Upload

After all taps and additional information is gathered, we show a view that the application is uploading the data to the server-side application.

4.2.2 Smartphone Sensors

Modern smartphones come with a variety of different sensors offering valuable services to it's users and enhancing many applications. The newest Apple iPhone to date, the iPhone 7, has a fingerprint sensor, a barometer, a three-axis gyroscope and accelerometer (MEMS), a proximity sensor and an ambient light sensor attached to it's main-board [12]. As we are going to predict finger taps on the iPhone screen, the only sensors that are effected by the force of the tap are the gyroscope and the accelerometer. Therefore, these will be outlined in the following sections.

Accelerometer

The accelerometer is a sensor module that measures the acceleration it encounters by either movement or gravity [21]. However, the acceleration caused by movement, the so-called inertial acceleration and the gravitational acceleration can not be distinguished by the sensor. This is due to Einstein's equivalence principal stating that the effects of gravity on an object are indistinguishable from the acceleration of the object's reference frame [34].

When the position of the device is fixed, as for example when it is placed on a table the accelerometer values would yield $a = \{a_x = 0, a_y = 0, a_z = -1\}$. This feature make it suitable for detecting device screen rotations. As the device flips from landscape to portrait orientation, the gravity is sensed by a different set of accelerometer axis [21].

The values of the acceleration are quantified in the SI unit metres per second per second (m/s^2). However, in engineering the acceleration is typically expressed in terms of the standard gravity (g).

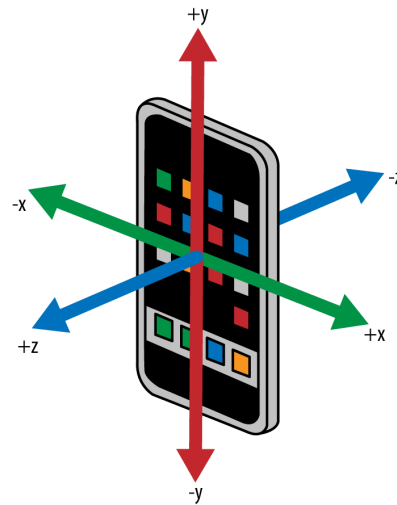


Figure 4.4: Apple iPhone with the corresponding axes of the accelerometer.

Gyroscope

As the accelerometer is suitable for detecting orientations, it lacks the ability to detect spin or more precise rotation movements. These spin movements are detected by the gyroscope sensor which is responsible for detecting and maintaining orientation [21].

A mechanical gyroscope typically composes of a spinning wheel which is set within three so-called gimbals. These gimbals enable the spinning wheel to be set in any orientation. Although the orientation does not remain fixed when the device is rotated, it changes in response to an external torque much lesser and in a different direction than it would be without the large angular momentum associated with it. Each gimbal translates to one of the three gyroscope outputs, namely *pitch*, *roll* and *yaw* (See [39]).

The gyroscope sensor within the MEMS², the chip deployed in the iPhone, is between 1 to 100 micrometers of size. When the gyroscope is rotated, a small resonating mass is shifted as the angular velocity changes. This movement is converted into very low-current electrical signals that can be amplified and read by a host system.

The values of the gyroscope are quantified in as rotations per seconds (RPS) or as degrees per second (deg/s).

²Microelectromechanical systems

Accessing sensor values

In order to access gyroscope and accelerometer Apple provides a high level API³ for accessing the device's sensors: **Core Motion**. Core Motion reports motion and environmental related data from sensors including accelerometers, gyroscopes, pedometers, magnetometers, and barometers in easy to use manner.

Sensor values can either be accessed as proceeded version including aggregations of the values and a raw version. For TapSensing, we make sure to record the raw values to avoid any form of bias. The update interval can be configured at ranges from 10Hz - 100Hz. Higher update-rates are possible but are not ensured to be processed in real-time by the device. For TapSensing, the update rate is configured with the highest (safe) value possible. This ensures that tap patterns are captured with an accurate resolution to make a later classification easier. The figure below is a code snippet depicting how sensor values are retrieved in the TapSensing application.

Figure 4.5: Swift code snippet displaying how to access sensor values with Core Motion.

Interoperability -JSON Scalability -Machine

4.3 Backend application

4.3.1 General requirements

Security

Consistency

4.3.2

4.3.3 Data model

4.3.4 User management & Security

4.3.5 Ensuring consistency

³An Application Program Interface is a set of rules and subroutines provided by an application system for the developer to use. Here is a link to the Core Motion API Documentation: <https://developer.apple.com/documentation/coremotion/>

5 Inferring Tap Locations

5.1 Data Preprocessing

5.2 Convolution Neural Network

5.3 Recurrent Neural Network

6 Method

6.1 Overview

6.2 Hypothesis

6.3 Experimental Setup

6.3.1 Subjects

6.3.2 Devices

6.3.3 Experiment Settings

Lab

Field

6.4 Analysis

7 Results

7.1 Hypothesis

7.2 Discussion

8 Conclusion

8.1 Further Outlook

List of Acronyms

3GPP	3rd Generation Partnership Project
AJAX	Asynchronous JavaScript and XML
API	Application Programming Interface
AS	Application Server
CSCF	Call Session Control Function
CSS	Cascading Stylesheets
DHTML	Dynamic HTML
DOM	Document Object Model
FOKUS	Fraunhofer Institut fuer offene Kommunikationssysteme
GUI	Graphical User Interface
GPS	Global Positioning System
GSM	Global System for Mobile Communication
HTML	Hypertext Markup Language
HSS	Home Subscriber Server
HTTP	Hypertext Transfer Protocol
I-CSCF	Interrogating-Call Session Control Function
IETF	Internet Engineering Task Force
IM	Instant Messaging
IMS	IP Multimedia Subsystem
IP	Internet Protocol
J2ME	Java Micro Edition
JDK	Java Developer Kit
JRE	Java Runtime Environment
JSON	JavaScript Object Notation
JSR	Java Specification Request
JVM	Java Virtual Machine
NGN	Next Generation Network
OMA	Open Mobile Alliance
P-CSCF	Proxy-Call Session Control Function
PDA	Personal Digital Assistant
PEEM	Policy Evaluation, Enforcement and Management
QoS	Quality of Service
S-CSCF	Serving-Call Session Control Function
SDK	Software Developer Kit
SDP	Session Description Protocol
SIP	Session Initiation Protocol
SMS	Short Message Service

SMSC	Short Message Service Center
SOAP	Simple Object Access Protocol
SWF	Shockwave Flash
SWT	Standard Widget Toolkit
TCP	Transmission Control Protocol
Telco API	Telecommunication API
TLS	Transport Layer Security
UMTS	Universal Mobile Telecommunication System
URI	Uniform Resource Identifier
VoIP	Voice over Internet Protocol
W3C	World Wide Web Consortium
WSDL	Web Service Description Language
XCAP	XML Configuration Access Protocol
XDMS	XML Document Management Server
XML	Extensible Markup Language

Annex

```
<?xml version="1.0" encoding="UTF-8"?>
<widget>
  <debug>off</debug>
  <window name="myWindow" title="Hello Widget" visible="true">
    <height>120</height>
    <width>320</width>
    <image src="Resources/orangebg.png">
      <name>orangebg</name>
      <hOffset>0</hOffset>
      <vOffset>0</vOffset>
    </image>
    <text>
      <name>myText</name>
      <data>Hello Widget</data>
      <color>#000000</color>
      <size>20</size>
      <vOffset>50</vOffset>
      <hOffset>120</hOffset>
    </text>
  </window>
</widget>
```

Listing 1: Sourcecode Listing

```

INVITE sip:bob@network.org SIP/2.0
Via: SIP/2.0/UDP 100.101.102.103:5060;branch=z9hG4bKmp17a
Max-Forwards: 70
To: Bob <sip:bob@network.org>
From: Alice <sip:alice@ims-network.org>;tag=42
Call-ID: 10@100.101.102.103
CSeq: 1 INVITE
Subject: How are you?
Contact: <sip:xyz@network.org>
Content-Type: application/sdp
Content-Length: 159
v=0
o=alice 2890844526 2890844526 IN IP4 100.101.102.103
s=Phone Call
t=0 0
c=IN IP4 100.101.102.103
m=audio 49170 RTP/AVP 0
a=rtpmap:0 PCMU/8000

SIP/2.0 200 OK
Via: SIP/2.0/UDP proxy.network.org:5060;branch=z9hG4bK83842.1
;received=100.101.102.105
Via: SIP/2.0/UDP 100.101.102.103:5060;branch=z9hG4bKmp17a
To: Bob <sip:bob@network.org>;tag=314159
From: Alice <sip:alice@network.org>;tag=42
Call-ID: 10@100.101.102.103
CSeq: 1 INVITE
Contact: <sip:foo@network.org>
Content-Type: application/sdp
Content-Length: 159
v=0
o=bob 2890844526 2890844526 IN IP4 200.201.202.203
s=Phone Call
c=IN IP4 200.201.202.203
t=0 0
m=audio 49172 RTP/AVP 0
a=rtpmap:0 PCMU/8000

```

Listing 2: SIP request and response packet[?]