

NTNU - NORGES TEKNISK-NATURVITENSKAPELIGE UNIVERSITET
Faculty of Engineering Science and Technology
Department of Civil and Transport Engineering
TBA4925 - Master Thesis

Optimizing the micro-tasking workflow and exploring it's usage potential within geospatial data

Anne Sofie Strand Erichsen
Trondheim, June 2017

DAIM page

Background

HEI

Task Description

The micro-tasking method is becoming more and more popular. Companies like Amazon develop micro-tasking web applications where people can earn money by doing micro-tasks for others. The method is used for tasks that involve both use of technology and a large number of people. By using the micro-tasking methodology, this thesis aims to study how people solves micro-tasks within geospatial data imports, which is a very complex and large process.

This study will have an emphasis on the data validation and conflict handling part of the import. These parts are complicated to do fully automatic through scripts. By varying the number of objects to solve at a time, adding rewards on some tasks, among other factors, the study will hopefully find a significant approach to prefer when using the micro-tasking method within geospatial data. What are the number of objects optimal within a task to get it completed as quickly as possible? Does the quality of the work vary between the different tasks given? Do amateurs manage to do the tasks? Do rewards have an impact on how the tasks are solved?

This thesis will also explore the micro-tasking methods usage potential within geospatial data. Can other organizations doing a process that needs humans to interfere take advantage of this method? An example is OpenStreetMap, who has taken good advantage of the method both in mapping and import projects.

Specific tasks:

- Study related literature
- Do a micro-tasking survey
- Examine how many elements are optimal when creating geospatial micro-tasks

Abstract

This paper propose a method for extracting buildings in satellite photos. The proposed network makes use of a digital surface model and multispectral satellite data. It

Sammendrag

Sammendrag på norsk

Preface

This paper is a master thesis written for the Department of Civil and Transport Engineering at the Norwegian University of Science and Technology (NTNU) in Trondheim, Norway. It is a part of the study program Engineering and ICT - Geomatics, and was written in the spring of 2017.

I would like to thank my supervisor Terje Midtbø for his help and feedback, and also Atle Frenvik Sveen for his support and help every time I needed it.

Trondheim, 2017-06-16?
Anne Sofie Strand Erichsen

Contents

Abstract	v
Sammendrag	vii
Preface	ix
1 Introduction	1
2 Background	3
2.1 Why do we need humans?	3
2.2 Human computation	4
2.3 Crowdsourcing	5
2.4 Micro-tasking	6
2.4.1 Micro-tasking platforms	6
2.4.2 Micro-tasking workforce	7
2.4.3 Micro-tasking usage	8
2.4.4 Building imports using micro-tasking	10
2.4.5 Challenges	10
3 Methodology and experiment	13
3.1 Experiment	13
3.1.1 Experiment questions	14
3.1.2 Experiment tasks	15
3.1.3 Building shapes	16
3.2 Web application	17
3.2.1 React application	17
3.2.2 Data acquisition	18
3.3 Pilot test	19
3.3.1 Execution of the pilot test	20
3.3.2 Results from the pilot test	20
3.3.3 Preliminary results	21
3.4 Sample Size	23
4 Result	25
4.1 Participants	25
4.2 Statistics theory	26
4.2.1 Normal testing	26
4.2.2 Hypothesis testing	28
4.3 Survey results	32
4.3.1 Gathered data	33
4.3.2 Normality tests	36
4.3.3 Levene's test of Equality of Variance	48
4.3.4 Hypothesis testing	49

CONTENTS

5	Discussion	63
6	Proposed sections	65
6.1	Future work	65
6.2	Usage potential	65
A	Appendices	67
A	Tets	69

List of Figures

2.1	Collective intelligence (Quinn and Bederson, 2011)	5
2.2	Micro-tasking (Michelucci and Dickinson, 2016)	6
2.3	Crisis map (Meier, 2014)	8
2.4	Swimming pools (Nikki, 2016)	9
3.1	Question one as it is displayed in the web application	14
3.2	Question two as it is displayed in the web application	14
3.3	Creating of footprint layers used in question one	16
3.4	Interface for the client application	18
3.5	Total time, all	22
3.6	Population vs. sample	23
4.1	The distribution of the task order in the analysed data	26
4.2	Skew (MedCalc Software bvba, 2017)	27
4.3	Kurtois (MedCalc Software bvba, 2017)	27
4.4	Histograms with normal distribution fit with samples containing total time to complete each task	37
4.5	Histograms with normal distribution fit after Box-Cox power transformation	38
4.6	Histograms with normal distribution fit with samples containing the number of correctly chosen elements in each task	39
4.7	Histogram with normal distribution fit - sample with total time per task	40
4.8	Histogram with normal distribution fit after Box-Cox transformation, sample with total time per task	41
4.9	Histogram with normal distribution fit showing samples with number of correct elements per task	41
4.10	Histogram with normal distribution fit after Box-Cox	42
4.11	Histograms with normal distribution fit with samples containing total time to complete each task	43
4.12	Histograms with normal distribution fit containing Box-Cox transformed data	43
4.13	Histogram with normal distribution fit showing samples with number of correct elements results for experienced participants	44
4.14	Histogram with normal distribution fit	45
4.15	Histogram with normal distribution fit after Box-Cox	45
4.16	Histogram with normal distribution fit	46
4.17	Histogram with normal distribution fit after Box-Cox transformation .	46
4.18	Sample 1 and 2 - mean (green dot) and standard deviation (blue line)	51
4.19	Sample 3 and 4 - mean (green dot) and standard deviation (blue line)	53
4.20	Sample 5, 6 and 7 - mean (green dot) and standard deviation (blue line)	55
4.21	Sample 8, 9 and 10 - mean (green dot) and standard deviation (blue line)	57

LIST OF FIGURES

4.22 Sample 11, 12 and 13 - mean (green dot) and standard deviation (blue line)	58
4.23 Sample 14, 15 and 16 - mean (green dot) and standard deviation (blue line)	59
4.24 Mean (green dot) and standard deviation (blue line) for sample 17, 18 and 19	60
4.25 Mean (green dot) and standard deviation (blue line) for sample 20, 21 and 22	61

List of Tables

4.1	Total time, all participants	33
4.2	Correct elements, all participants	34
4.3	Total time, divided into task 1, 2 and 3	34
4.4	Correct elements, divided into task 1, task 2 and task 3	34
4.5	Total time, task and experienved divided	35
4.6	Correct elements, task and experienved divided	35
4.7	Total time, inexperienced per task	36
4.8	Correct elements, inexperienced per task	36
4.9	Summary, normality tests	47
4.10	Summary, Levene's tests	48
4.11	Summary, hypothesis tests	61

1 | Introduction

To the authors best knowledge, little research has been done on micro-tasking geospatial data. This thesis aim is to study how well geospatial data tasks is solved through micro-tasking. The quality of the completed tasks when doing micro-tasks it is important. In this study the quality is measured through the number of correctly chosen elements in each task. The resulting data will distinguish between experienced and inexperienced participants. A micro-task should be small enough so that all individuals can complete the task, independent on their background and experience.

Salk et al. (2016) looked at how local knowledge and professional background, impact the volunteer performance in a Land-Cover classification task. The paper concluded that there was no difference in how well the participants did, and their background.

A task can traditionally be divided by time, place, person, object, and skill [(Meier, 2013b), p. 13]. A task can be created by identifying the time it will require, the place where it must be done, the people who need to do the task, the object on which the work is done and finally, the skill needed for the task. Today we have technology that can create and move tasks around based on the four first categories. Technology can allocate tasks based on the deadline and time the task requires. It can also establish communication between any team of people dependent on which people the task requires. One thing that technology can't do is change the skill of individual workers, though it can only connect people with different skills to work on the same tasks [(Meier, 2013b), p. 14]. Crowdsourcing moves beyond this and looks at the skills of individual workers, the problem that needs to be solved and combines the best skills of workers to solve the problem. The division of labor by skill has more economic impact than the other four categories [(Meier, 2013b), p. 15]. Crowdsourcing is a way of refactoring work in a way that exploits the worker's flexibility and gets the right skills to the right part of the problem. To get the right skills to the right part of the problem it needs to be partitioned into smaller parts. Having smaller parts will make it easier to distribute the problem. The distribution can be done through micro-tasking, also called "smart crowdsourcing" by Patrick Meier (Meier, 2013a).

This thesis aim is to study if micro-tasking can successfully be expanded to involving maps and geospatial data. The OpenStreetMap community has used the method some time, and the usage so far can be evaluated as successful. This thesis also aims to find out if inexperienced individuals also manage to solve tasks on maps that involves geospatial data. The study also aims to determine if the number of elements in each micro-task has an impact on how well individuals solve the micro-tasks. The quality of the work, the number of correctly solved tasks and time, is measured. The thesis uses a survey hosted through a web-application to gather participant data. The data is then used to answer this thesis hypothesizes. The next chapter will give a thorough introduction to micro-tasking and hopefully make it clearer what this thesis aim is. Chapter 3 will explain the survey and chapter 4 will contain the statistics,

1. INTRODUCTION

both hypothesis, theory, and results.

Albert Einstein illustrates this perfectly: “Computers are incredibly fast, accurate, but stupid. Humans are incredibly slow, inaccurate, but brilliant. Together they may be powerful beyond imagination” (Holzinger, 2013).

2 | Background

Creating and maintaining real-world knowledge bases in a classical work environment demands a high cost, and is a cost that is often unnecessary [(Meier, 2013b), p. 134]. Alternative approaches are to rely on the knowledge of open crowds, volunteer contributions, or services like micro-tasking platforms where there are people ready to work on the tasks given to them [(Meier, 2013b), p. 134].

Today, geospatial data is more available than ever. Governments are releasing more and more data and the OpenStreetMap database is still growing. While general data availability is increasing, the quality of the data is not necessarily perfect and manual pre-processing is often necessary before using it (Difallah et al., 2015). Pre-processing of the data can require much time and high costs. By exploiting both machines and people through the appropriate platform and approach, the cost can decrease and the quality increase. The author will argue in this chapter that combining machines and people is often a better and faster solution than a fully-automatic or fully-manual approach and implementing such an approach into a micro-tasking platform can be a good solution.

2.1 Why do we need humans?

Machine learning gives computers the ability to learn without being explicitly programmed. It involves computer intelligence, but the computers do not know the answers up front (Stanford University, 2017). Machine-learning algorithms have enormous problems when contextual information is missing. Without a pre-set of rules, a machine has trouble solving the problem. Machines do not have creativity, which is required to answer complex problems (Holzinger et al., 2016). According to the company "Mighty AI," humans cannot be removed from Artificial Intelligence training loops. They believe that humans will continue to play a crucial role in creating training data for the algorithms.

It is suggested by Biewald (2015) and Oppenheimer (2017) that machine learning accuracy should follow the Pareto 80:20 principle. Getting 80 % accuracy can be reasonably easy to accomplish, but the last 20 % should be handled by human input (Biewald, 2015) (Oppenheimer, 2017). Human input can be to label the original training dataset or help correct inaccurate predictions outputted from the algorithm (Oppenheimer, 2017). The machine learning company "developmentSEED" use a micro-tasking solution to clean their machine learning output data. They are using humans to get a faster, more accurate output data. "developmentSEED" developed Skynet Scrubber, a GUI web application solution to get human input quicker and easier (Their algorithm is called Skynet). In their blog, Derek Lieu writes: "Skynet gets more capable every day, but the output is still not perfect [...] We built Skynet Scrub so we could start using Skynet data sooner".

2. BACKGROUND

Holzinger (2016) claim that most people from the machine learning community are concentrating on *automatic* machine learning by bringing the humans out of the process. When humans are out of the loop, the training data sets can be uncertain and incomplete, and the resulting algorithm can be questionable (Holzinger, 2016). By bringing humans back in the process, especially in domains where the data sets are questionable, for instance in the health domain, one enables what neither a human or a computer can do on their own (Holzinger, 2016). It is today possible to build hybrid human-machine systems that combine both the scalability of computers and the yet unmatched cognitive abilities of the human brain (Difallah et al., 2016). "Computers are bad at finding patterns unless we have a well-understood problem" quote Stephen Cohen, co-founder of Palantir Technologies, only humans can understand and frame a new problem. Palantir Technologies believe in augmenting human intelligence, not replacing it. As Holzinger (2013) say, "[...] the problem-solving knowledge is located in the human mind and - not in machines, " and this is something we must acknowledge according to Holzinger.

2.2 Human computation

Human computing is, at its most general level, computation performed by human beings and a human computation system contains both humans and computers working together to solve difficult problems (Schulze et al., 2012). The author argues that utilizing the human processing power is still important. Humans are necessary even though our computers are becoming more and more complex. Traditional approaches to solving problems are to focus on improving the software, but as the reader will see in this paper, a solution that uses humans cleverly by exploiting the human brain's cognitive abilities can create much faster and better results than software. One of the pioneers of crowdsourcing, Luis von Ahn, wanted to find a cheap and efficient way to label images (von Ahn, 2008). The solution was to exploit the use of a game-like approach in a non-game context to motivate individuals to label the pictures through a game. This approach is called gamification (Huotari and Hamari, 2017). The game was called "The ESP game" and solved the problem of labeling images with words. Most images do not have a proper caption associated with them, and this makes it difficult to create search engines for images. A fast and cheap method of labeling images is by using humans cleverly, and humans can very easily see if the picture contains, i.e., a dog or a cat. Through "The ESP game" humans where labeling images without even knowing it, they only played a fun game. Within a few months, the game collected more than 40 million image labels (von Ahn, 2008), and they did not even have to pay them doing it. Human computation is one of the major areas where the gamification approach has been employed (Morschheuser et al., 2016). Each human performs a small part of a massive computation task.

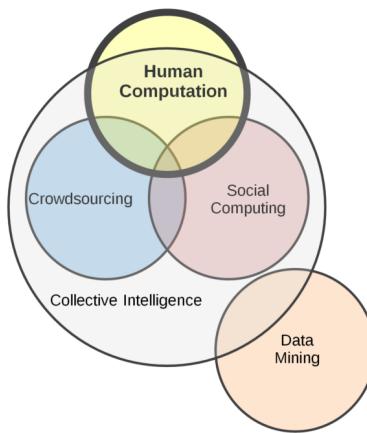


Figure 2.1: Collective intelligence (Quinn and Bederson, 2011)

Human computation, a term introduced by Luis von Ahn, refers to, according to Quinn and Bederson (2011), a distributed system that combine the strengths of humans and computers to accomplish tasks that neither can do alone. To make human computation in crowdsourcing compelling one needs to know how the results can be optimally acquired from humans and how the results can be integrated into productive environments without having to change established workflows and practices [(Meier, 2013b), p. 134]. Gamification can be one solution on how to make human computation in crowdsourcing effective (Wang et al., 2017). The author will argue that micro-tasking can be another solution for effective crowdsourcing using humans cognitive abilities.

2.3 Crowdsourcing

The first time the term "crowdsourcing" appeared was in Wired magazine article by Jeff Howe (Howe, 2006). Whereas human computing (2.2) replaces computers with humans, crowdsourcing replaces traditional human workers with members of the public (Quinn and Bederson, 2011). EYeka (2015) state that 85 % of the top global brands use crowdsourcing for various purposes. Crowdsourcing is an increasingly important concept (Salk et al., 2016) and has become a widespread approach to dealing with machine-based computations where we leverage the human intelligence (Gadiraju et al., 2015).

When the scope of a crowdsourced project is explicitly geographical, it is often called *volunteered geographical information* (VGI). According to Salk et al. (2016), the best known VGI project is OpenStreetMap (OSM). OSM is an open-source mapping project, where volunteers contribute with their local knowledge and mapping abilities.

2.4 Micro-tasking

The simplest type of tasks are called micro-tasks and are illustrated in figure 2.2. Micro-tasks should not require any special training, and a task should be completed within a couple of minutes (Ipeirotis and G., 2010). Problems that are suitable for solving through micro-tasking are those that are easy to distribute into many simple tasks, which can be completed in parallel in a relatively short period of time, without requiring specific skills (Sarasua et al., 2012). Research has also demonstrated that micro-tasking is effective for far more complex problems when using sophisticated workflow management techniques. Micro-tasking can then be applied to a broader range of challenges like: (1) completing surveys, (2) translating text between two languages, (3) matching pictures of people, (4) summarizing text (Bernstein et al., 2015).

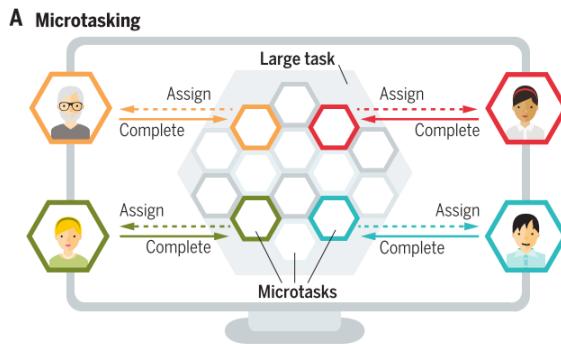


Figure 2.2: Micro-tasking (Michelucci and Dickinson, 2016)

2.4.1 Micro-tasking platforms

2.4.1.1 Amazon's Mechanical Turk

Amazon's Mechanical Turk (MTurk) is a very popular micro-tasking platform created in 2005, but still in use today (Difallah et al., 2016). MTurk acts as an online labor marketplace (Sarasua et al., 2012). It provides the infrastructure, connectivity and payment mechanisms so that hundreds of thousands of people can perform micro-tasks on the Internet and get paid per completed task. MTurk is used for many different tasks that are easier for people than computers. It contains simple tasks such as labeling or segmenting images or tagging content, to more complex tasks such as translating or even editing text (Franklin et al., 2011). In the marketplace, employers are known as requesters, and they post tasks, called *human intelligence tasks* (HIT's). The HIT's are then picked up by online users, *crowd workers*, who complete the tasks in exchange for a small payment (a few cents per HIT) (Ipeirotis and G., 2010).

2.4.1.2 Tasking manager

The Tasking Manager tool is OpenStreetMap's micro-tasking platform. It was created in the aftermath of the Haiti earthquake in 2010 (Palen et al., 2015). The tool is used to coordinate satellite image tracing projects and sorts the area covered by the satellite image into grids so that multiple people can map the same area at the same time. Each person works at one grid each, this way they avoid mapping the same areas. Organizing the areas into grids is a very effective approach to coordinate a mapping job. The tasking manager is mainly used by the *Humanitarian OpenStreetMap Team* (HOT). This platform does not have a rewarding system or a gamification approach. It is solely based on volunteer contributors. This platform shows that it is not necessary to have a game or rewards for a successive platform.

There are also other tools in OpenStreetMap. Tofix etcetc.

2.4.1.3 CrowdFlower

CrowdFlower is a company that wants to help businesses take advantage of crowdsourcing and human computation. They act as an intermediary for these companies (Quinn and Bederson, 2011). CrowdFlower receives tasks from businesses wanting to crowdsource their work or problems. CrowdFlower operates with a variety of services to get connected with workers (i.e., MTurk) (Quinn and Bederson, 2011).

What's special with CorwdFlower is their close ties with AI technology and a crowd-sourced workforce. Their costumers are allowed to perform tasks with algorithms and machine learning, but bring in human judgment when they are not confident in the technology, and the human work can make the algorithms smarter (Ha, 2016). The founder of CrowdFlower says that "self-driving cars have gotten pretty good at recognizing many of the objects they encounter on the street, [...] (but) they can still struggle with tricky things like "a person in a Halloween costume dressed as a stationary object, or a pole with a person painted on it," which is where CrowdFlower comes in." (Ha, 2016).

2.4.2 Micro-tasking workforce

It is said that crowdsourcing is radically changing the nature of work Deng et al. (2016). Traditional workers are restricted to offices and arranged office hours. With crowdsourcing, through for instance micro-tasking platforms, the workers can choose when to work, and even better: which jobs to perform. This appears very attractive, but is it only on the surface?

According to Deng et al. (2016), evidence indicates that crowdsourcing is radically changing people's perspectives on how to manage their work-life balance. Compared to "traditional" work tasks, the micro-tasks are simple and fast to finish (within a couple of minutes). The worker is also often compensated with tiny rewards every time they complete a micro-task, which is motivating.

2. BACKGROUND

Individuals who perform micro-tasks for micropayment is called *crowd workers* by (Deng et al., 2016). A study done on workers in the micro-tasking platform MTurk (section 2.4.1.1), says that the workers are representative for the general Internet user population, but are generally younger and have lower incomes and smaller families (Ipeirotis and G., 2010).

2.4.3 Micro-tasking usage

Micro-tasking and human computation have close ties. In the "Handbook of Human Computation," micro-tasking is strongly present in the *Human Computation for Disaster Response* chapter [(Meier, 2013b), p. 95-105], as well as in several other parts of the book. In the disaster response chapter, the authors give an overview of how human computation methods, such as paid micro-tasks, could be used to help in major disasters. In 2012, Philippines was struck by a typhoon called Ruby, devastating large regions. Through CrowdFlower the workers collected over 20 000 tweets related to the typhoon and identified the tweets containing links to either photos or video footage from the damaged areas. The photos or videos in the relevant tweets were tagged and geo-tagged by volunteers if they portrayed evidence of damage. Within 12 hours a dataset of 100 georeferenced images and videos were collected. It resulted in a very detailed crisis map shown in figure 2.3. This map was the first official crisis-map based solely on social media content [(Meier, 2013b), p. 101]. In the aftermath of this crisis, an algorithm was developed to automatically detect tweets that link to photos and videos, which freed more time for the volunteers to georeference and tag more images and videos portraying evidence of damage (Meier, 2014).

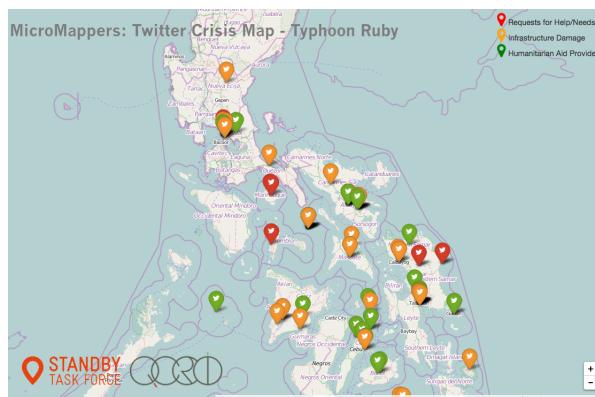


Figure 2.3: Typhoon Ruby Crisis map (Meier, 2014)

Micro-task crowdsourcing refers to a problem-solving model in which a problem or task is outsourced to a distributed group of people by splitting the task or problem into smaller sub-tasks or sub-problems. The sub-tasks is solved by multiple workers independently, often in return for a reward (Sarasua et al., 2012). Thanks to micro-tasking platforms as MTurk, it is possible to build a hybrid human-machine system

that combines the scalability of computers with the yet unmatched cognitive abilities of the human brain (Difallah et al., 2016). Gadiraju et al. (2015) findings when analyzing data from MTurk, indicate rapid growth in micro-task crowdsourcing. With the establishment of micro-task crowdsourcing platforms as MTurk and CrowdFlower, micro-tasking is much more accessible. Micro-tasking practitioners are actively turning towards paid crowdsourcing to solve data-centric tasks that require human input (Gadiraju et al., 2015). Most cases of micro-tasking combine human computation abilities with crowdsourcing.

Human computation systems often use crowdsourcing platforms to recruit workers (Schulze et al., 2012). Companies developing machine-learning algorithms has seen an advantage in combining human computation abilities and crowdsourcing with fast machine learning algorithms. An example is the team "Tomnod"¹ who did a project in Australia where they combined human computation, crowdsourcing, and machine learning to locate swimming pools (Kostas, 2016). The machine learning algorithm classified polygons where there was likely to be a swimming pool inside. The crowd who participated in finding the swimming pools were only presented with the classified polygons, which minimized the search area and then also the time required finishing the job (Kostas, 2016). Examples of classified polygons are shown in figure 2.4. The "Tomnod" team divided the task into smaller tasks, they micro-tasked the work. The resulting dataset was then used to train a swimming pool detecting convolutional neural network (Nikki, 2016).



Figure 2.4: Classified polygons created by the algorithm (Nikki, 2016)

Most cases of micro-tasking usage exploit the large volume capabilities machines have and the cognitive capabilities of humans (Difallah et al., 2016). One of the advantages of micro-tasking platforms like MTurk, "Tomnod" and CrowdFlower, mentioned by Meier (2013b) (p. 99), is the built-in quality control mechanisms that ensure a relatively high quality of output data. They set a review constraint, for instance in a project where they tagged satellite imagery of Somalia each unique image was re-

¹Tomnod is a team of volunteers who work together to identify important objects and interesting places in satellite images; www.tomnod.com

2. BACKGROUND

viewed by at least three different volunteers and only when all three agreed on type and location it was approved.

2.4.4 Building imports using micro-tasking

In OpenStreetMap, at least two building imports were successfully completed using micro-tasking. The tasking manager platform (2.4.1.2) was used to organize the import.

New York Public Library uses micro-tasking to train computers to recognize building shapes and other data on digitized insurance atlases. The micro tasks are used to check the computer's work and also to capture information the computer missed. Individuals contributing checks and fixes building footprints drawn by the computer. The individuals also enter addresses and classify the building footprints using colors. To ensure accuracy the same footprints are shown to several people. At least three different individuals check the same footprint, and 75% or more must agree on the footprint for the answer to be approved.

2.4.5 Challenges

Getting enough people to use the micro-tasking platforms is crucial for its success. Most of the platforms mentioned in this chapter give payments to the workers. Another option is to make the platform as a game, which is also shown in this chapter. Creating a micro-tasking platform without payments or gamification factors the page is likely to have a short life, even though the tasking manager, supported by HOT, is an exception to this rule.

A problem when combining machines and humans is that machines can do their operations in real-time, while humans are unpredictable, they can come and go as they wish. This creates a gap where the micro-tasking platforms cannot guarantee on the task completion time (Difallah et al., 2016).

The human computation abilities can also be overestimated. During the classification of swimming pools in Australia, the Tomnod team faced some unexpected challenges. As described in section 2.4, they used the crowd to classify if a polygon contained a swimming pool or not, an algorithm had pointed out the polygons first. When reviewing a random sample from the result, they found an indication that 26% of polygons that contained a pool were identified as not containing pools by the crowd (Kostas, 2016). Further studies also showed that the guilty part was the crowd, the algorithm had correctly detected polygons containing pools. In a case where the algorithm was 85% confident that the polygon contained a pool, only one voted 'yes', six voted 'no', this polygon does not contain a pool. The solution was to combine the human verdict with the machine's prediction. This example shows that it is important to use the right combination of humans and machines. Tasks that at first seem simple to do for humans, may be more challenging than expected. Basic object detection using machine learning perform very well when used together with human operations.

It is important that operations added to a micro-tasking platform consider the talents and limitations of human workers (Franklin et al., 2011) and this is what this thesis try to examine. What are the limitations of human workers when dealing with maps and geospatial data. It has been shown that crowds can be "programmed" to execute classical algorithms such as Quicksort, but such use of available resources is neither performant nor cost-efficient (Franklin et al., 2011).

New software developed by researchers at Facebook can score 97.25 percent on the same challenge, regardless of variations in lighting or whether the person in the picture is directly facing the camera."

3 | Methodology and experiment

There is little research on how inexperienced individuals solve micro-tasks when the tasks involve map interaction. To the authors best knowledge, little, if any, research has been conducted on micro-tasks involving geospatial interpretation and analysis. An experiment will be carried out in this thesis, and the results will be used to get more knowledge about micro-tasks involving geospatial tasks.

Gadiraju et al. (2015) categorize the top-level crowdsourced tasks, after analyzing platforms as MTurk and CrowdFlower. It resulted in six classes, three relevant classes within geospatial data is *1) Verification and validation*, *2) Interpretation and analysis* and *3) Content creation*. There are examples of all three task classes in geospatial crowdsourcing. During imports of large datasets into OpenStreetMap, crowdsourcing is used to validate the new data. In humanitarian OpenStreetMap, they map areas during a crisis to support the help organizations through crowdsourcing, creating valuable content to the workers in the field (ref xx). In a machine learning process, they are starting to use micro-tasks to both validate the created data and also create data sets to train the algorithms. Shown in chapter 2. *Interpretation and analysis* tasks rely on the individual to use their interpretation skills during task completion. The experiment conducted in this thesis use Gadiraju et al. (2015) three classes, with emphasis on the *Interpretation and analysis* class, to develop the questions given to the participants. This section will introduce the experiment developed. The experiment generates the data for the analyses, so it is important to implemented and executed the experiment correctly.

3.1 Experiment

In general, questionnaires are used to generate quantitative data, which is later used to calculate statistical information [(Kitchin and Tate, 2000), p. 48]. The experiment is used to generate data and to answer hypothesis involving geospatial micro-tasks. An experiment containing three tasks answering was developed to be able to explain the hypothesis introduced in chapter one.

The three tasks will vary the number of elements the participant has to use to answer the questions. Each task will contain the same two questions. The questions represent two separate micro-tasks involving geospatial data. The participant will always answer the two questions on six elements, but the tasks vary how many elements to handled at the same time. The variation of the number of elements in the tasks is to hopefully find out if or how much the number of elements in a micro-task effect how well individuals solves the task.

3. METHODOLOGY AND EXPERIMENT

3.1.1 Experiment questions

The experiment asks two questions, which represents two different geospatial micro-tasks. As mentioned in the introduction in this chapter, the *Interpretation and analysis* task class, listed by Gadiraju et al. (2015), is emphasized and gives the context to both questions. The first question asks the participant to click on the color that fits the shape of the marked building(s) on the map best. Question one is an analyzing task. The participant is given two footprint layers on each building and needs to determine which of the footprints that fit the building shape shown on the base map best. Question one displayed on the web application is presented in figure 3.1. In this example, the participants have three buildings to select.

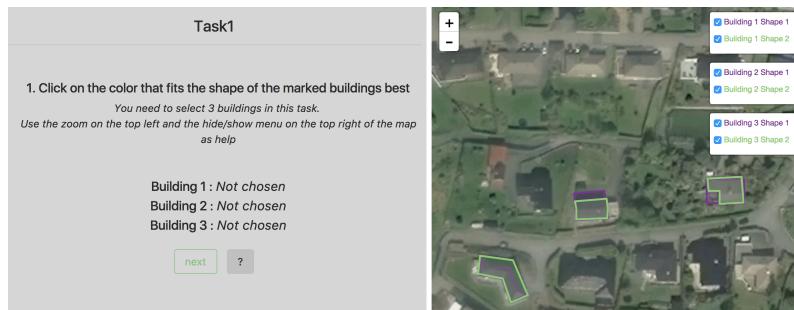


Figure 3.1: Question one as it is displayed in the web application

The second question asks the participant to select the 1/3/6 most informative row(s) that describes random buildings best. Question two is an interpretation task where the participant needs to interpret the information written in the table to decide which row(s) gives the most informative information about an arbitrary building. In figure 3.2 question two is displayed in the web application. In this example, the participant has to choose three rows from the table.

Task1			
2. Select the 3 most informative rows that describes random buildings best			
<i>Each row represents a new building. Think that the information should be informative for everyone, independent of education, background etc.</i>			
Choose	Info 1	Info 2	Info 3
<input type="checkbox"/>	Height: 9 m	Gnr: 33	Bnr: 169
<input type="checkbox"/>	Country: Norway	City: Bergen	Address: Hammarslandgrenda 66
<input type="checkbox"/>	Validation date: 20160816	Registered: Yes	Area: 1015,9
<input type="checkbox"/>	Building type: Detached house	Building levels: 3	Building material: Brick and wood
<input type="checkbox"/>	Address: haMarslaNgrenda	Municipality: Bergen	Country: Unknown
<input type="checkbox"/>	Source: Photogrammetric data capture	Building: House	Amenity: Place of residence

Figure 3.2: Question two as it is displayed in the web application

3.1.2 Experiment tasks

The experiment consists of three different tasks, in addition to a training task. When determining the number of elements in the three tasks the author decided to base this on cognitive load theory. Cognitive load theory refers to the total amount of mental effort being used in the working memory. Working memory is determined by the number of information elements that need to be processed simultaneously within a certain amount of time (Barrouillet et al., 2007). A heavy cognitive load can have negative effects on task completion. The cognitive load that is imposed by a task is much higher for beginners than for more advanced students (Leppink et al., 2014).

The survey will contain three tasks, each task contains six elements, but the task varies how many elements the participant need to answer at the same time. One task will serve the participant with one and one element, demanding the smallest cognitive load. The other task will serve the participant with three and three elements at the same time. This number is just under the limit of how much information humans can process. The last task will serve the participant with all six elements at the same time. This number exceeds the human capacity when processing information according to Leppink et al. (2014).

It is stated that the working memory has a limited capacity of seven plus or minus two elements (or chunks) of information when merely holding information and even fewer (ca four) when processing information (Leppink et al., 2014). By choosing three elements in one task and six elements in the other task, this paper can determine if the theories about the limited capacity of the human brain also apply to micro-tasks involving geospatial data. One task will only contain one element as a minimum cognitive load task. The three tasks can help answer how many elements a human can process at the same time without impacting the quality of the result. The goal is to determine a preferred number of elements to include in a micro-task. This information can be useful when developing micro-tasks to achieve adequate task progress as well as accurate results.

The “magical” number of 4 has been demonstrated to limit much of human information processing (Mandler, 2013). It is said that polygon comparison demand medium cognitive load (Kiefer et al., 2016), which is what the participants do in the first question in the experiment. Kiefer et al. (2016) argues that high cognitive load may lead to less efficient map reading and spatial orientation, as well as decreased spatial learning. Since polygon comparison demand medium cognitive load, question one should at least not be too demanding on the one- and the three elements task. A worry is that the inexperienced participants will have a bigger struggle than the experienced participants. The extraneous cognitive load imposed high for the inexperienced when solving problems because their lack of prior knowledge of how to solve that type of problem forces them to resort to weak problem-solving strategies (Leppink et al., 2014). By dividing the participants into experienced and inexperienced categories, the results from the experiment can help determine if geospatial micro-tasks are too demanding on inexperienced individuals.

3. METHODOLOGY AND EXPERIMENT

3.1.3 Determining the building footprints used in question one

According to Fan et al. (2014), there was over 77 million buildings in the OpenStreetMap (OSM) database in 2013. A study of the geometries of building footprints in the city Munich reveal a huge diversity in the geometries (Fan et al., 2014), and this is probably not the only city with this kind of diversity. The Fan et al. (2014) paper used four criterion's, completeness, semantic accuracy, position accuracy and shape accuracy to evaluate the quality of the building footprints in OSM. In the creation of the two footprints representing the buildings used in question one, the quality criterion's shape- and position accuracy was emphasized. The goal is to create shapes that match realistic cases that occur for instance in OSM.

Shape accuracy evaluates how well the layer matches the building in an aerial image. Fan et al. (2014) mentions two main reasons to why building footprints are simplified in OSM. The first reason is the difficulties following building details when looking from a bird's eye view. The second reason is the limited resolution of the Bing aerial image used in OSM during digitalization. In question one two footprints is drawn with one of them matching the building shape better than the other. The participant has to use an aerial image to determine which layer fits the building best. This will test if the participants manage to make correct shape judgments by only using an aerial image as a reference.

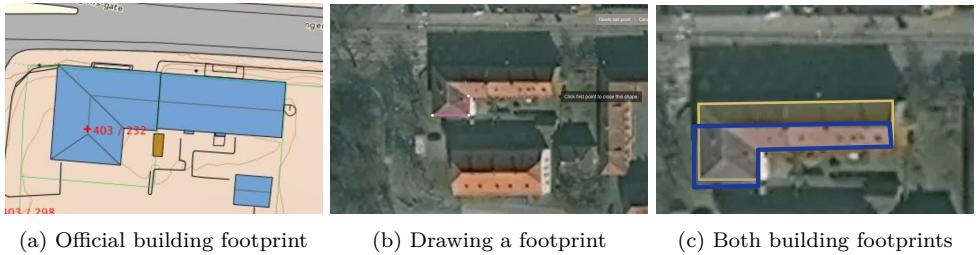


Figure 3.3: Creating of footprint layers used in question one

Position accuracy evaluates how well the coordinate value of a building relates to the reality on the ground. Fan et al. (2014) tested the accuracy of buildings in OSM and concluded with a mean offset of 4.13 m. The low positional accuracy of OSM building footprints data is caused by the limited resolution of Bing map images. By combining shape- and position accuracy in some of the cases used in question one this study can also determine if participants manage to evaluate both factors. In this study, the participants do not have available information about what the true ground coordinates are. Therefore position accuracy will be examined by shifting one of the layers. The correct positional accuracy will be at the building in the aerial image.

3.2 Web application

This thesis used an online web-based survey to conduct the experiment. An online survey avoids the cost and effort of printing, distributing, and collecting paper forms. Many people prefer to answer a brief survey displayed on a screen instead of filling in and returning a printed form (Ben and Plaisant, 2009). The participants do not have to share the same geographic location as the researcher.

An online web environment also makes it easier to use interactive maps. An interactive map is necessary to answer question one. It is not possible to have an experiment involving interactive maps on a piece of paper. The web is the obvious way of implementing interactive maps. Making it online, available via URL, makes the distribution faster and easier.

A common web programming language is JavaScript, together with the library React¹ creates the client side of this thesis application. The client communicates with a server that fetches the task elements from and saves the task results in a PostGIS database. The server is written in Python with the framework Django². The PostGIS database contains the task elements, and also the task results gathered from the participants.

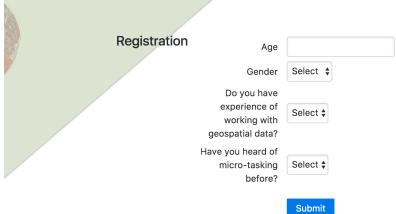
3.2.1 React application

The React application was created to serve the experiment to all participants. It contains all the steps of the experiment. First, the participant registers, giving information about age, gender and answers yes or no on the following two questions: 1) "Do you have experience of working with geospatial data?", 2) "Have you heard of micro-tasking before?". Next, the participant is given an introduction page with a detailed introduction video describing how to answer the two questions, on how to interact with the map and building layers. A training task comes after the instruction video, in the training task the participant solves both questions just like the normal tasks, only it contains different building footprints and only two elements in each task to not replicate the other tasks. After the training task, the participant starts with the experiment containing the three tasks, with a short survey after each task. The survey asks the participant to rate the task difficulty between one and five, and if the participant tried his/hers best or was interrupted during the task. The participant can also write a comment. Interfaces for the two questions is shown in figure 3.1 and 3.2 in section 3.1.1. The registration form and survey interface is shown in figure 3.4.

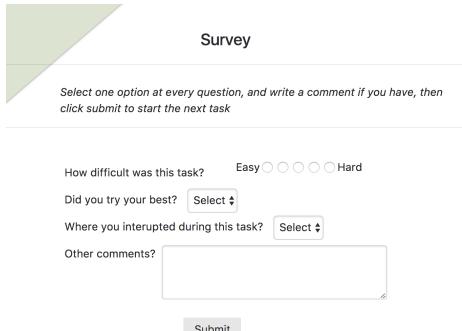
¹React is a open-source JavaScript library for building user interfaces. In React, the displayed data can change without reloading the page. It's main goal is to be fast, simple and scalable.

²Django is a open-source web framework written in Python. It's primary goal is to ease the creation of complex, database-driven websites.

3. METHODOLOGY AND EXPERIMENT



(a) Registration site



Survey

Select one option at every question, and write a comment if you have, then click submit to start the next task

How difficult was this task? Easy Hard

Did you try your best?

Where you interrupted during this task?

Other comments?

Submit

Figure 3.4: Interface for the client application

Multiple measurements to ensure random and independent observations were implemented during the development of the React application. The main measurements to ensure random, independent observations was:

- Random order on the three tasks
- Random buildings in the tasks
- Random color on the building footprints
- Random order on which the building footprints was drawn on the map
- Random order on where the information rows was written in the table

When distributing an online experiment, the result can be inconsistent since the researcher is not present to either control the participant or the environment surrounding the participant. The researcher implemented functions securing the completeness of the data. The buttons navigating to the next question and task was disabled until enough building footprints or rows were selected. The participant could not submit their task-result before answering both questions correctly. The submit button on the register form and after each task survey only submitted the answers if all fields were answered. If a field missed a message appeared asking the participant to fill out the entire form before submitting. All registered results were complete in the database thanks to these measurements.

3.2.2 Data acquisition

In this study, the independent variables are: 1) experienced or inexperienced participant, 2) number of elements in the task, 3) age and 4) gender of the participant. Independent variables are factors we think might influence the results of the study [Kitchin and Tate, 2000], p. 49]. When a participant registers at the start of the survey the independent variables is generated. Participants who answer yes to the

question "Do you have experience of working with geospatial data?" are registered as experienced. The independent variables are believed to influence the dependent variables. Dependent variables are factors the study is interested in explaining [(Kitchin and Tate, 2000), p. 49]. In this study, the dependent variables are: 1) time spent on each question and task, 2) the number of correctly chosen elements in each question and task, and 3) how difficult the participant thought the task was.

The React application generates and saves the task results. In both questions, the time and number of correct elements are registered. The React application has a timer that registers time elapsed on both questions and adds the time measurements together to get the total time spent on the task. The time only reflects how long the participants spent solving the two questions. Time spent loading new layers, moving to the next question, etc., is not included. Which building footprints and rows the participant selects is also registered, and before saving the result counts the number of correctly chosen elements in each question and then adds the number together to get the total number of correctly chosen elements in the task. A participant can maximum have twelve correct elements, six from each question. Total time and number of correct elements are the two primary dependent variables, and they create the basis of the statistical analyses together with the independent variables mentioned above. Participants task results are saved after each task and contain the following information:

- Task number (which task)
- Task order number (which order)
- Time spent on question one
- Time spent on question two
- Total time spent on the task
- Correct elements in question one
- Correct elements in question two
- Total correct elements in the task

After each task the participants answers a short survey. The survey asked the participant about the difficulty of the task, if the participant tried it's best or was interrupted during the task and the participant can also write a comment. This information was used to remove task results where the participants was interrupted. The difficulty question can be used to determine if one of the three tasks is preferred by the participants.

3.3 Pilot test

Testing the experiment before actual use is highly recommended (Ben and Plaisant, 2009). A pilot test provides an opportunity to validate the wording of the tasks.

3. METHODOLOGY AND EXPERIMENT

It also helps understand the time necessary for completing the survey, which should be communicated to the participants (Schade, 2015). The pilot test was carried out on a small sample of users. Results from the pilot test were in this thesis used to make improvements to the actual survey, to the react application and to find errors or weaknesses in the database models.

After the pilot test, the usability was measured. Usability in this thesis was measured with the *System Usability Scale*(SUS) because it gives a subjective measure of usability. The *System Usability Scale* questionnaire consists of ten statements where the participants rate their agreement on a five-point scale (Ben and Plaisant, 2009). Subjective measure of usability is usually obtained through the use of questionnaire and scales (Brooke, 1996). SUS was developed to be quick and straightforward, but also reliable enough to be able to compare performance changes between versions (Brooke, 1996). It is also easy to administer the participants through the usability test, and it can be used on small sample sizes and still give reliable results (Affairs, 2013).

The usability is important to measure. If the participants do not understand how the web application works, they will probably not do the survey since they then have to invest time in understanding what to do. It is also important to get enough participants to do the whole survey and not quit halfway in frustration of not understanding it properly. The *System Usability Scale* can effectively differentiate between usable and unusable systems (Affairs, 2013).

3.3.1 Execution of the pilot test

The pilot test was conducted with a total of eight participants, five experienced and three non-experienced participants aged from 22 to 64 years. It started with brief information about this study and the experiment. They were told to talk out loud during the test, no help or guidance was given to the participants. The author observed the participants while they conducted the survey. The author took notes and watched if the participants understood the questions and tasks correctly. After the survey a *System Usability Scale* questionnaire was answered by the participants. In the end, the participants were asked to give general feedback on the web application. The SUS score and the feedback were then used to determine the usability of the React application and to determine which improvements to be done.

3.3.2 Results from the pilot test

The average SUS score was 84.64 out of 100. Anything above 68 is considered above average (Affairs, 2013). When adding the SUS score to an adjective rating will an score of 85.5 or higher be described as excellent (Bangor et al., 2009). A score of 84.64 is then described as good/excellent. This result gives a strong indication that the React application is user-friendly.

All participants thought that the instruction video was confusing. It was short, the instructions went too fast, and it missed voice descriptions. The instructions needed

major improvements, an important discovery. The purpose of the video is to give the participant an introduction to how to answer the two questions. It should include important instructions, particularly useful for participants not used to working with interactive maps.

Overall feedback on question one was that it is hard to understand which building was which because of missing labels, and also to know when a building footprint was selected or not. The lack of labels on the buildings was done on purpose to get the task as much as possible realistic. The process of selecting the best fitting building footprint needed improvements. It had to be clearer that one chosen by clicking on the layer on the map, not by using the layer control as some thought. This part was added to the movie with voice description, describing in detail how a footprint was selected. Improvements to the design of the question one page was made by adding the selected footprint's color to the text telling the participant which layer they had chosen, giving a better visual feedback.

Another feedback from one of the participants was that both question views had too much information and long sentences. The participant advised to shorten the sentences and to move some of the information to the video.

The test data was used to examine some of the hypothesis to find errors or weaknesses in the database model. The data was extracted from the PostGIS database and saved in CSV files. Some preliminary results can be seen in section 3.3.3. There were a few errors and weaknesses found during the statistical tests. Changes to the database models were necessary, and the implemented changes are listed under:

1. Add foreign key from TaskResult model to TaskSurvey model
2. Added four other fields in TaskResult model
 - Total correct elements
 - Task order
 - Task number
 - List of correctly chosen building numbers in both questions
3. The difficulty field in Tasksurvey model was changed from Char field to Integer field

3.3.3 Preliminary results

The pilot-test data was not normally distributed, and doing statistical analysis on data from eight participants didn't seem relevant. The data was mapped in char plots, visualizing some trends.

The two oldest participants spent almost twice as much time on the test than the younger. Maybe it was too much cognitive load on them. Learning a new application and at the same time understanding how to do the survey and answer the questions given to them. One of them were experienced and the other inexperienced, so this is

3. METHODOLOGY AND EXPERIMENT

a surprising result. Figure 3.5 show the task results from all participants ordered by age. There are three entries per participant, so three and three bars are results from the same participant. Task 1 represents the task with one elements, task 2 the task with three elements and task 3 the task with six elements. It is clear to see that the 54 year old participant's total time on task 2 is dramatically higher than the rest.

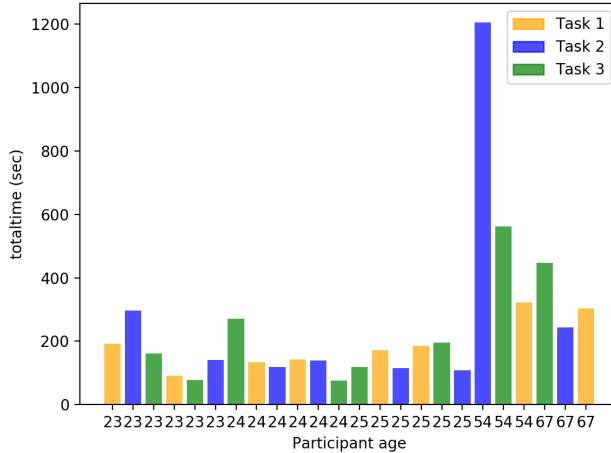


Figure 3.5: Total time - all participants ordered by age

The average time spent on the survey was 18 minutes. The two oldest participants used on average 33 minutes, while the rest of the participants spent on average 13 minutes to complete the survey.

In the pilot test, the same building footprints (question one) and information rows (question two) were used in all three tasks. At the end of the pilot test, the author asked the participants if they remembered the buildings and meta information from the previous tasks. $\frac{7}{8}$ answered yes on the question. This information was valuable. If every participant does a better job at the last task, because they remember the elements from earlier tasks, the result will not be as useful. Reading the data in figure 3.5, $\frac{6}{8}$ participants spent less time on the last task, even though the task order varied. This almost matches the number of participants who remembered the previous elements in the last task. This finding made the author create three different task element groups. The three element groups were randomly assigned to each task, to avoid this decision to influence the result. The risk of participants remembering previous task elements disappeared.

3.4 Determining the sample size

The sample size is influenced by various factors, including the purpose of the study, population size, the risk of selecting a "bad" sample and the allowable sampling error (Israel, 1992).

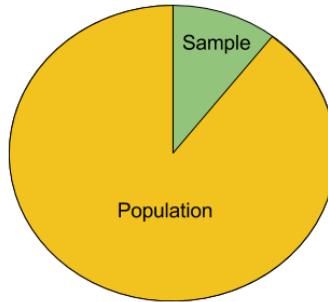


Figure 3.6: Population vs. sample

A sample is a collection of observations and is the subset of a population, illustrated in figure 3.6. The population size in this survey is not easily determined. A population is the collection of individuals of a particular type (Walpole et al., 2012). All individuals with access to a computer and the internet interested in contributing to micro-tasks can be one description of the population. It is important that the sampled population and the target population is similar to one another.

There are three possible ways of determining the sample size in this study. The first option is to use a sample size from a similar study. The risk is to repeat errors that were made in determining the sample for another study. The second option is to rely on published tables, depending on precision, confidence levels, and variability. According to Israel (1992) table 1, an accuracy of 0.05, confidence interval of 95% and a size of population greater than 100'000, the necessary sample size is 400. If the accuracy is changed to 0.1, the sample size necessary increases to 100 (Israel, 1992). The numbers found in the table must reflect the number of obtained responses. The last approach is to use formulas to calculate the sample size. The formulas require the standard deviation and how much variance to expect in the response (Smith, 2013)(Israel, 1992). Israel (1992) mentions that the table gives a useful guide for determining the sample size and that formulas are used if the study has a different combination of precision and confidence. This study will use the table result since the combinations match this study.

It is important to mention that the quality of the sample is as important as the size. The more variable the sampled data is, the larger the sample size is required (Israel, 1992). It is also desirable to choose a random sample, which means that the observations are made independently and random. The main purpose of using a random sample is to obtain correct information about the unknown population parameters (Walpole et al., 2012).

4 | Result

This chapter will present the result of the experiment conducted. It will summarize who the participants were, the gathered data and then present the statistical results from the analysis performed on the gathered data. The researcher used eight days to collect enough participants completing the experiment. Only 38% of the registered participants completed the experiment. All task results saved in the database was valid and could be used in the analyses. The analyses were calculated in Python, using statistical packages as SciPy, Numpy, and Panda. The data was extracted from the database using Django Queryset and saved in CSV files. On the authors GitHub, in the repository *thesis-statisticmethods*, the implemented statistical methods are available¹.

4.1 Participants

Benefits of using an online based experiment are that it has potential to reach a huge number of people, the problem is how to reach out to the people to make them aware of the existence of the application. Using mailing threads and sites the researcher had available was this studies solution. The participant's contribution to the survey was contacted through email and the organization Geoforum's website and Facebook page. All students at *Civil and engineering* was emailed, as well as a mailing list reaching out to the Norwegian OpenStreetMap community. The web application was also published at the Geoforum website (www.geoforum.no) and Facebook page. Geoforum is a Norwegian association for individuals and companies working in the field of geomatics. The participants receiving the email or looking at the Geoforum site could click on the web application URL and access the experiment from there.

There was in total 461 task results in the database after the gathering period, where almost all participants contributed with three task results (the training task was removed). 402 participants registered on the website during the data gathering period and only 38% of the registered participants completed the survey. 152 participants completed all three tasks. Results from participants not completing all three tasks are also included in the dataset. Including these result is not a problem. The three tasks were given to each participant in a random order, and the tasks are also independent, containing different buildings and rows.

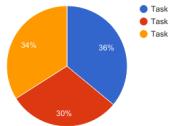
The mean participant age was 31.5 years and the median 25 years. The youngest participant was 19 and the oldest 58 years. 33% of the participants was female and 66% male. The average male was 33.5 years old and the average female 31.2 years old. 19% of the participants that completed the survey said they had heard of micro-tasking before. 53% of the participants stated that they had experience of working with

¹<https://github.com/annesofie/thesis-statisticmethods>

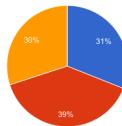
4. RESULT

geospatial data. The distribution between experienced and inexperienced participants was approximately even. A very pleasing distribution.

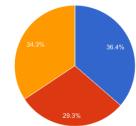
Random and independent observations are essential, and a random task order was used to ensure this. Analyzing the task results and the order the tasks was presented in gives a pleasing result. The distribution of how many times the three tasks (Task 1, 2 and 3) was first, second and last task in the analyzed data is approximately evenly distributed. This is shown in figure 4.1.



(a) First task



(b) Second task



(c) Third task

Figure 4.1: The distribution of the task order in the analysed data

4.2 Statistics theory

This section will give an introduction to the statistics used in this thesis. The thesis will examine the data with parametric methods but also with non-parametric methods if the assumption of a normally distributed samples fails. A nonparametric method is much more efficient than the parametric procedure when the set of data used in the test deviates significantly from the normal distribution (Walpole et al., 2012). There are also some disadvantages using nonparametric methods. The methods will be less efficient, and to achieve the same power as the corresponding parametric method a larger sample size is required. If parametric and nonparametric tests are both valid on the same set of data, the parametric test should be used (Walpole et al., 2012).

4.2.1 Normal testing

The sampling distribution of a statistic depend on the distribution of the population, the size of the samples, and the method of choosing the samples (Walpole et al., 2012). Sampling distribution describes the variability of sample averages around the population mean μ . All parametric statistics assumes normally distributed, independent observations. Parametric tests are preferred in statistics because it got more statistical power than nonparametric tests (Frost, 2015). The power of a test is the probability of correctly rejecting a false null hypothesis, which in this case is the ability to detect if the sample comes from a non-normal distribution. To determine if a sample is normally distributed there exists both visual methods and normality tests to assess the samples normality. A visual inspection of the sample's distribution is usually unreliable and does not guarantee that the distribution is normal (Pearson

et al., 2006). Presenting the data visually gives the reader an opportunity to judge the distribution themselves. In this thesis histograms are used to visualize the data for normality.

Normality tests compare the scores in the sample to a normally distributed set of scores with the same mean and standard deviation (Ghasemi and Zahediasl, 2012). There are multiple normality tests, and deciding which test to use is not easy. This study needs a test that doesn't require every value to be unique, a test that can handle ties (identical observations). The survey used to collect the samples in this study do not guarantee unique values.

The D'Agostino-Pearson omnibus test stand out as the best choice. This test first computes the skewness, see figure 4.2, and kurtois, see figure 4.3, to quantify how far from the normal distribution the sample is from the terms of assymetry and shape. Then it calculates how far each of these values differs from the value expected with a normal distribution (Pearson et al., 2006). It works well even if all values are not unique (Motulsky, 2013). The test also works well on both short- and long-tailed distributions (Yap and Sim, 2011).

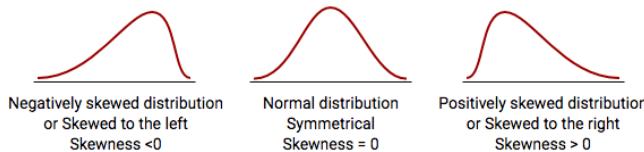


Figure 4.2: Skew (MedCalc Software bvba, 2017)

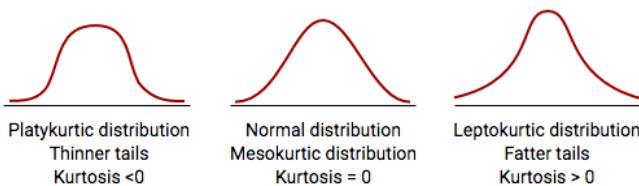


Figure 4.3: Kurtois (MedCalc Software bvba, 2017)

The D'Agostino-Pearson test uses the following hypothesis:

$$H_0: \text{The data follows the normal distribution}$$

$$H_A: \text{The data do not follow the normal distribution}$$

For small sample sizes, normality tests have little power to reject the null hypothesis, therefore small sample sizes most often pass normality tests. For large sample sizes, significant results would be derived even in the case of a small deviation from normality (Pearson et al., 2006). When the null hypothesis cannot be rejected, then there are

4. RESULT

two possible cases. First case is to accept the null hypothesis or the second case is that the sample size is not large enough to either accept or reject the null hypothesis (The Pennsylvania State University, 2017). An acceptance of the null hypothesis implies that the evidence was insufficient, the result does not necessarily accept H_0 , but fails to reject H_0 (Walpole et al., 2012).

4.2.2 Hypothesis testing

The null- and alternative hypothesis are statements regarding a difference or an effect that occur in the population of the study. The alternative hypothesis (H_a) usually represents the question to be answered or the theory to be tested, while the null hypothesis (H_0) nullifies or opposes H_a (Walpole et al., 2012). The sample collected in the study is used to examine which statement is most likely (technically it is testing the evidence against the null hypothesis). When the hypothesis is identified, both null and alternative, the next step is to find evidence and develop a strategy for or against the null hypothesis (Lund Research Ltd, 2013a).

The next step, after the hypothesis is identified, is to determine the level of statistical significance, often expressed as the *p-value*. A statistical test will result in the probability (*the p-value*) of observing your sample results given that the null hypothesis is true. A significance level widely used in academic research is 0.05 or 0.01 (Walpole et al., 2012).

The result should not be reported as "significantly different," but instead report it as "statistically significantly different." This is because the statistical decision as to whether the result is significant should not be based solely on the statistical test. To indicate to readers that the result is a statistical one, include statistically in the conclusion sentence (Lund Research Ltd, 2013c).

4.2.2.1 Two-sample t-test

When estimating the difference between two means a two-sample t-test is used (Walpole et al., 2012). A two sample test assumes two independent, random samples from distributions with means $[\mu_1, \mu_2]$ and variances $[\sigma_1^2, \sigma_2^2]$. The hypothesis tested on two means can be written as

$$\begin{aligned} H_0: \mu_1 - \mu_2 &= 0 \text{ or } \mu_1 = \mu_2 \\ H_A: \mu_1 - \mu_2 &\neq 0 \text{ or } \mu_1 - \mu_2 > 0 \text{ or } \mu_1 > \mu_2 \end{aligned}$$

The two-sample t-test is used to estimate if differences between two means are significant. In a two-sample, two-sided, t-test ($\mu_1 - \mu_2 \neq 0$) the null hypothesis is rejected when [(Walpole et al., 2012), p. 345]:

$$|T| > t_{\frac{\alpha}{2}, v} \quad (4.1)$$

In a two-sample, one-sided, t-test the null hypothesis is rejected when [(Walpole et al., 2012), p. 350]:

$$T > t_{\frac{\alpha}{2}, v} \quad (4.2)$$

$$T < -t_{\frac{\alpha}{2}, v} \quad (4.3)$$

Equation 4.2 is used on one sample test where the alternative test is to check if the mean is greater than zero ($\mu_1 - \mu_2 > 0$), and the 4.3 equation is used on hypothesis where the test is to check if the mean is lower than zero ($\mu_1 - \mu_2 < 0$). T is the calculated statistical value and t is the critical value with the given significance level (α) and degree of freedom (v). The critical value is found in the table of Critical values for t-distribution.

Before doing tests on the two means, the Levene's Test is used to test if the samples are from populations with equal variances. It tests the hypothesis:

$$\begin{aligned} H_0: & \text{Input samples are from populations with equal variances} \\ H_A: & \text{Input samples are from populations that do not have equal variances} \end{aligned}$$

If we can assume equal variances in the two samples and the samples are normal distributed, a two-sampled t-test may be used.

Relevant hypothesis in this study that can be tested with a two-sampled t-test (if the conditions mentioned above are valid) is listed under.

Hypothesis - Two sample t-test

$$\begin{aligned} H_0: & \text{Experienced and inexperienced spent the same amount of time on the tasks} \\ H_A: & \text{Total task time differs between them} \end{aligned}$$

$$\begin{aligned} H_0: & \text{Experienced do not finish the tasks more quickly than inexperienced} \\ H_A: & \text{Experienced participants finish the tasks faster} \end{aligned}$$

$$\begin{aligned} H_0: & \text{Total number of correct elements between experienced and inexperienced are equal} \\ H_A: & \text{There is a difference in number of correct elements between them} \end{aligned}$$

$$\begin{aligned} H_0: & \text{Experienced no not have more total correct elements then inexperienced} \\ H_A: & \text{Experienced participants have a higher number of correct elements} \end{aligned}$$

4. RESULT

Before solving the hypothesis the conditions needs to be testet. More on this later.

4.2.2.2 Analysis-of-Variance

Analysis-of-Variance (*ANOVA*) is according to Walpole et al. (2012) a very common procedure used for testing population means. Where a two sample t-test are restricted to consider no more than two population parameters, *ANOVA* can test multiple population parameters. A part of the goal of *ANOVA* is to determine if the differences among the means of two or more samples are what we would expect due to random variation alone, or due to variation beyond merely random effects. *ANOVA* assumes normally distributed, independent, samples with equal variance. The equal variance assumption will be tested with Levene's Test also mentioned in subsection 4.2.2.1.

One-way *ANOVA* tests the null hypothesis that two or more groups have the same population mean given that the mean is measured on the same factor or variable in all groups(Lund Research Ltd, 2013c). The hypothesis test can be written like this:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$
$$H_A: \text{At least two of the means are different}$$

μ equals the group mean and k represents the number of groups. It is important to check that each group are normally distributed (Lund Research Ltd, 2013c). The weakness of one-way *ANOVA* is that it cannot tell which specific groups were significantly different from each other if H_0 is rejected. To be able to determine which group a *post hoc test* is used. The null hypothesis is accepted if:

$$F - \text{statistic} < f_{\alpha, v_1, v_2}(\text{critical value}) \quad (4.4)$$

If the alternative hypothesis is accepted a *post hoc test* is used. A post hoc test makes paired comparisons to determine which groups differs. This thesis will use Tukey's test to determine which groups means are significantly different [(Walpole et al., 2012), p.526].

In a one-way *ANOVA* test there should be one dependent variable and minimum three independent groups, which is an relevant approach considering the data produced from this thesis survey. There are at least two dependent variables in the survey data, task time and number of correctly chosen elements. The survey result can be divided into three groups, one element task, three elements task and six elements task. Each entry in the sample should only be assigned to one group. Relevant hypothesis from the study that can be used in an one-way *ANOVA* analysis is shown under.

Hypothesis - One-way ANOVA

H_0 : Mean task time is not different between the three tasks
 H_A : Mean task time is different between at least two of the tasks
Variable = time, group = tasks

H_0 : Total number of correct elements between the three tasks are equal
 H_A : Total number of correct elements between at least two of the tasks are not equal
Variable = Number of correct elements, group = tasks

The hypothesis written above will be tested in section 4.3.4.3.

4.2.2.3 Wilcoxon Rank-Sum test

The Wilcoxon Rank-Sum test is an appropriate alternative to the two-sample t-test (see subsection 4.2.2.1) when the normality assumptions do not hold, but the samples are still independent and have a continuous distribution (Walpole et al., 2012). Since this method is nonparametric (or distribution-free) it does not require the assumption of normality.

The hypothesis for Wilcoxon Rank-Sum Test is:

$$H_0: \tilde{\mu}_1 = \tilde{\mu}_2 \\ H_A: \tilde{\mu}_1 > \tilde{\mu}_2 \text{ or } \tilde{\mu}_1 < \tilde{\mu}_2 \text{ or } \tilde{\mu}_1 \neq \tilde{\mu}_2$$

The alternative hypothesis depends on what the test should determine. If the sample with mean $\tilde{\mu}_1$ is greater than, smaller than or unequal to the sample with mean $\tilde{\mu}_2$. First select a random sample from each population with means $\tilde{\mu}_1$ and $\tilde{\mu}_2$. If the sample sizes are different, let n_1 be the number of observations in the smallest sample and n_2 for the largest sample. Then $\tilde{\mu}_1$ will be the mean for the smallest sample. If there are ties (identical observations) in the sample a Mann-Whitney U test is preferred (The Scipy community, 2017).

4.2.2.4 Mann-Whitey U test

The Mann-Whitney U test is used to compare differences between two independent groups. This test can be used to conclude whether two populations differ. It can for instance test if there are differences in medians between groups (Lund Research Ltd, 2013b). In contrast to the t-test, it compares the median scores of two samples instead of the mean score. The test is non-parametric and can therefore be used on samples that are not normally distributed. The test assumes that the samples come from populations with equal variances. When comparing two sample medians the two independent variables (i.e experienced and inexperienced participants) has to have a

4. RESULT

similar shape. It can test the hypothesis:

$$H_0: \text{The two populations are equal}$$
$$H_A: \text{The two populations are not equal}$$

The null hypothesis is rejected if (LaMorte, 2017):

$$U \leq \text{critical value} \quad (4.5)$$

The critical value is found in the table of Critical Values for U and depends on the sample sizes, n_1 and n_2 , and the significant level α . U is the statistical value calculated.

4.2.2.5 Kruskal-Wallis test

The Kruskal-Wallis test is a nonparametric alternative to one-way *ANOVA* (see subsection 4.2.2.2) (Walpole et al., 2012). This test should be used if the assumption of normal distribution failed. As mentioned in this sections introduction, a nonparametric method does not assume normality. This test is an generalization of the rank-sum test when there are more than 2 samples.

Kruskal-Wallis is used to test equality of means in one-way *ANOVA*, so the hypothesis for the Kruskal-Wallis test is:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$
$$H_A: \text{Minimum two of the } \mu_k \text{'s are different}$$

Here μ_k is the rank mean for the group k. As in Wilcoxon Rank-Sum test (subsection 4.2.2.3), the number of observations in the smallest sample is assigned to n_1 , the second smallest to n_2 and the largest sample is assigned to n_k .

The null hypothesis is accepted if:

$$H - \text{value} < h_\alpha \quad (4.6)$$

4.3 Survey results

- All participants ordered by age, excluded by task 4
- All results in one task, ordered by age
- Average time per micro-task
- Is there a difference in task number 1, 2, 3? time and correct
Can use it to explain the data

4.3.1 Gathered data

The gathered data is analyzed on the two dependent variables: 1) total time used to complete each task and 2) number of correctly chosen elements per task. Both variables sum the participants time spent and correct elements on question one and question two together. The gathered data is also divided by the two independent variable pairs: 1) experienced- and inexperienced participants and 2) the three survey tasks. Combining the dependent- and independent variables create the foundation of the 22 different samples in this thesis. The samples are listed in the tables in this section. Sample mean \bar{x} , sample median, standard deviation of \bar{x} , standard error ($\frac{\text{standard deviation}}{\sqrt{\text{samplesize}}}$) of \bar{x} , minimum in sample and maximum in sample is listed in the tables. The sample number is also written in the tables and is used in the analysis of the data, to easier distinguish which sample is used in the different analyses. In all the samples, results from the training task and from participants that were disturbed during the survey is removed.

The gathered data are in the first subsection (4.3.1.1) divided into experienced and inexperienced. In subsection 4.3.1.2 the data is divided by the three tasks, containing all participants. Section 4.3.1.3 and 4.3.1.4 divides the data by the three tasks and also in experienced and inexperienced participants.

4.3.1.1 All, experienced and inexperienced participants

Table 4.1 and 4.2 are samples containing task results from all, experienced and inexperienced participants. The result is divided into the two dependent variables, total time and the number of correctly chosen elements.

Total time per task (seconds) Sample number	All	Experienced 1	Inexperienced 2
Number of observations	429	229	200
Sample mean \bar{x}	170.32	177.65	161.94
Sample median	155.0	158.0	154.0
Standard deviation of \bar{x}	82.19	88.24	73.99
Standard error of \bar{x}	3.98	5.83	5.23
Minimum in sample	38.00	52.00	38.00
Maximum in sample	657.00	657.00	529.00

Table 4.1: Total time spent on each task

4. RESULT

<i>Correct elements per task</i> Sample number	All 3	Experienced 3	Inexperienced 4
Number of observations	429	229	200
Sample mean \bar{x}	9.82	9.81	9.83
Sample median	10.0	10.0	10.0
Standard deviation of \bar{x}	1.52	1.53	1.51
Standard error of \bar{x}	0.07	0.10	0.11
Minimum in sample	4.00	5.00	4.00
Maximum in sample	12.00	12.00	12.00

Table 4.2: Number of correctly chosen elements per task

4.3.1.2 All participants, divided by task

In table 4.3 and 4.4 the task results is divided into the three different tasks. Task 1 is the task that served the participants with one and one elements. Task 2 is the task that served the participants with three and three elements, and task 3 gave all six elements at the same time. Table 4.3 contains the total time variable, and 4.4 the number of correct elements variable.

<i>Total time per task (seconds)</i> Sample number	Task 1 5	Task 2 6	Task 3 7
Number of observations	146	142	141
Sample mean \bar{x}	166.38	172.25	172.48
Sample median	150.0	155.5	157.0
Standard deviation of \bar{x}	84.57	84.21	77.95
Standard error of \bar{x}	7.00	7.07	6.56
Minimum in sample	47	50	38
Maximum in sample	657	492	529

Table 4.3: Total time divided into task 1, task 2 and task 3

<i>Correct elements per task</i> Sample number	Task 1 8	Task 2 9	Task 3 10
Number of observations	146	142	141
Sample mean \bar{x}	10.19	9.71	9.55
Sample median	11.0	10.0	10.0
Standard deviation of \bar{x}	1.43	1.53	1.52
Standard error of \bar{x}	0.12	0.13	0.13
Minimum in sample	5.00	5.00	4.00
Maximum in sample	12.00	12.00	12.00

Table 4.4: Number of correctly chosen elements divided into task 1, task 2 and task 3

4.3.1.3 Experienced participants, divided by task

In the tables in this section, only task results from experienced participants are included, and the result is also divided into the three survey tasks. Table 4.5 is data gathered about experienced participants total time per task. Table 4.6 is data gathered about experienced participants number of correctly chosen elements per task.

<i>Total time per task</i>	Task 1	Task 2	Task 3
Sample number	11	12	13
Number of observations	77	80	77
Sample mean \bar{x}	173.04	176.70	181.06
Sample median	158.0	156.0	165.0
Standard deviation of \bar{x}	96.76	86.13	79.70
Standard error of \bar{x}	11.03	9.63	9.08
Minimum in sample	57.00	52.00	53.00
Maximum in sample	657.00	492.00	463.00
Sample number	11	12	13

Table 4.5: Experienced total time per task, divided by task

<i>Correct elements per task</i>	Task 1	Task 2	Task 3
Sample number	14	15	16
Number of observations	77	80	77
Sample mean \bar{x}	10.29	9.66	9.48
Sample median	11.0	10.0	10.0
Standard deviation of \bar{x}	1.32	1.64	1.47
Standard error of \bar{x}	0.15	0.18	0.17
Minimum in sample	7.00	5.00	5.00
Maximum in sample	12.00	12.00	12.00

Table 4.6: Experienced number of correct elements per task, divided by task

4.3.1.4 Inexperienced participants, divided by task

In this section, the task results from only inexperienced participants are included, and the result is also divided into the three survey tasks. Table 4.7 is the total time variable and 4.8 the number of correctly chosen elements.

4. RESULT

<i>Total time per task (seconds)</i>	Task 1	Task 2	Task 3
Sample number	17	18	19
Number of observations	71	64	65
Sample mean \bar{x}	158.30	165.69	162.23
Sample median	148.0	154.5	154.0
Standard deviation of \bar{x}	67.57	80.93	74.53
Standard error of \bar{x}	8.02	10.12	9.24
Minimum in sample	47.00	50.00	38.00
Maximum in sample	487.00	455.00	529.00

Table 4.7: Inexperienced total time per task, divided by task

<i>Correct elements per task</i>	Task 1	Task 2	Task 3
Sample number	20	21	22
Number of observations	71	64	65
Sample mean \bar{x}	10.07	9.78	9.61
Sample median	10.0	10.0	10.0
Standard deviation of \bar{x}	1.54	1.38	1.57
Standard error of \bar{x}	0.18	0.17	0.19
Minimum in sample	5.00	6.00	4.00
Maximum in sample	12.00	12.00	12.00

Table 4.8: Inexperienced number of correct elements per task, divided by task

4.3.2 Normality tests

The samples need to be tested if they follow a normal distribution. This test is important for determining if a parametric or a non-parametric test should be used when including the different samples in the hypothesis tests. The section about normal testing (4.2.1) concluded that the D'Agostino and Pearson normality test should be used in this thesis. A visual interpretation of histograms will also be a part of the normality test. D'Agostino-Pearson uses the following hypothesis:

$$H_0: \text{The data follows the normal distribution}$$

$$H_A: \text{The data do not follow the normal distribution}$$

4.3.2.1 Experienced and inexperienced participants, total time variable

This section will test if sample 1 and 2 (table 4.1) follow a normal distribution. Sample 1 and 2 will be used to determine if there are a significant difference in time spent on the tasks between experienced and inexperienced participants.

A visual interpretation of histogram 4.4a and 4.4b gives an indication that sample 1 and 2 are not normally distributed. Both histograms are positively skewed (figure 4.2). Samples involving time measurements are rarely normally distributed. This is because the samples will always be skewed since it is impossible to have negative time and there will always be a limit to how fast a participant can finish the task.

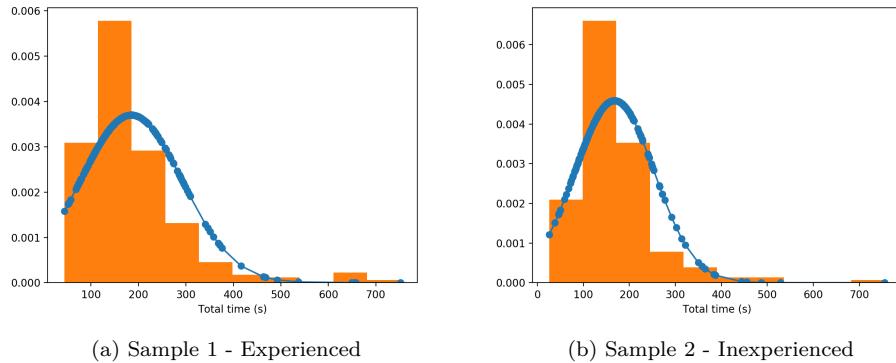


Figure 4.4: Histograms with normal distribution fit with samples containing total time to complete each task

A D'Agostino and Pearson normality test confirmed the visual interpretation with a significance level of 5%. Both samples obtained p-values lower than the significance level, and the null hypothesis is rejected. Sample 1 and 2 are not normally distributed with a confidence level of 95%.

D'Agostino and Pearson normality test
Significance level: 5%

Sample 1

P-value: $3.874 * 10^{-22}$

The p-value is lower than the significance level (0.05), the null hypothesis is rejected and H_A accepted.

Sample 2

P-value: $2.574 * 10^{-21}$

The p-value is lower than the significance level (0.05), the null hypothesis is rejected and H_A accepted.

In both sample 1 and 2, the p-value was significantly lower than the significance level of 5%. Data transformations are commonly used tools to improve normality of a sample's

4. RESULT

distributions, but there are many types of data transformations. Osborne (2010) claim that almost all tests, even non-parametric tests, benefit from improving the normality of the samples, especially when the normality test is significantly denied. Typical traditional transformations are square root, inverse or converting to logarithmic scales (Osborne, 2010).

A Box-Cox power transformation (Box-Cox) is used in this thesis. This transformation can only be used on positive data and the data gathered in this thesis will never be below zero. Box-Cox takes the idea of having a range of power transformations (square root $x^{\frac{1}{2}}$, inverse x^{-1} , etc.) available to improve the effectiveness of normalizing and variance equalizing for both positively- and negatively-skewed variables (Osborne, 2010). This transformation will always use the appropriate conversion to be maximally effective in moving each sampled data towards normality. This is the reason why this thesis will use the Box-Cox transformation.

Sample 1 and 2 after a Box-Cox is shown in histogram 4.5a and 4.5b. A visual inspection gives a good indication that the transformed data is normally distributed. The skewness looks approximately zero.

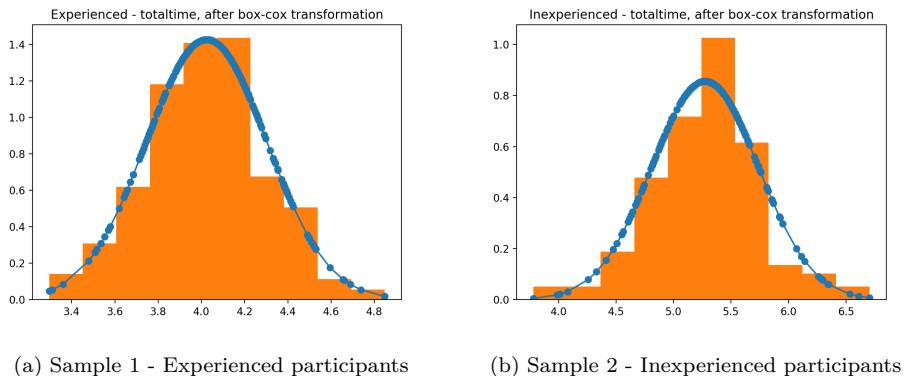


Figure 4.5: Histograms with normal distribution fit after Box-Cox power transformation

The transformed data is then applied to the D'Agostino and Pearson test. This test confirms the visual analysis, both sample 1 and sample 2 are normally distributed after the Box-Cox with a confidence level of 95%. The calculated p-value is larger than the significance level of 5%.

D'Agostino and Pearson normality test

(After Box-Cox transformation)

Significance level: 5%

Sample 1: Experienced, total time per task

P-value: 0.849

The p-value is higher than the significance level (0.05), the null hypothesis is accepted.

Sample 2: Inexperienced, total time per task

P-value: 0.0623

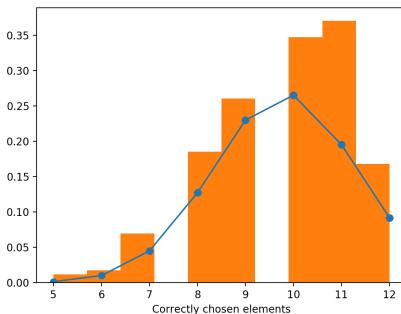
The p-value is higher than the significance level (0.05), the null hypothesis is accepted.

The assumption that sample 1 and sample 2 are normally distributed is now accepted and the transformed data can be used in parametric methods.

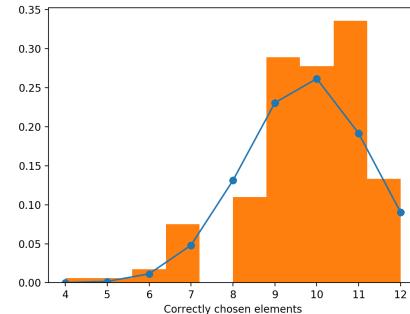
4.3.2.2 Experienced and inexperienced participants, number of correctly chosen elements variable

This section will test if sample 3 and 4 (table 4.2) are normally distributed. Sample 3 and 4 will be used to test if there are any difference in the number of correctly chosen elements per task between experienced and inexperienced participants.

A visual inspection of the samples histogram 4.6a and 4.6b gives a good indication that sample 3 and 4 are not normally distributed. Both are negatively skewed (figure 4.2).



(a) Sample 3 - Experienced



(b) Sample 4 - Inexperienced

Figure 4.6: Histograms with normal distribution fit with samples containing the number of correctly chosen elements in each task

4. RESULT

D'Agostino and Pearson normality test confirm our visual analysis. Both samples accept the alternative hypothesis with p-values (0.00443, 0.00013) lower than the significant level 0.05. The null hypothesis is rejected and H_A accepted for sample 3 and 4.

Sample 3 and 4 is Box-Cox power transformed because the null hypothesis was rejected. After the transformation, a new D'Agostino and Pearson normality test was performed. Both samples also failed this test. Sample 3 and 4 are not normally distributed and need to be tested with non-parametric methods.

4.3.2.3 All participants divided by task, total time variable

In this section sample 5, 6 and 7 (table 4.3) is normality tested. These samples will be used to test whether there is a significant difference between the three tasks when looking at the total time variable.

A visual analysis of the three histograms in figure 4.7a, 4.7b and 4.7c show a positive skewness, just like the histograms in figure 4.4. This gives an indication that the three samples are not normally distributed.

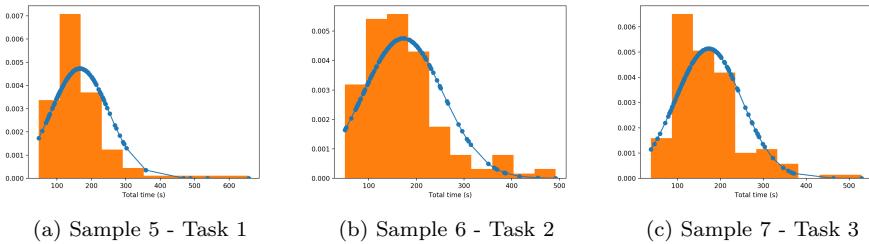


Figure 4.7: Histogram with normal distribution fit - sample with total time per task

The D'Agostino and Pearson normality test agreed with the visual analysis. Obtained p-values for all three samples ($2.39 * 10^{-24}$, $2.57 * 10^{-9}$ and $1.71 * 10^{-11}$) are smaller than the significance level 0.05, and the null hypothesis is rejected. The samples are not normally distributed with a significant level of 5%.

Because the null hypothesis was rejected, the samples are Box-Cox power transformed. Histograms of each sample after the transformation is shown in figure 4.8a, 4.8b and 4.8c. A visual analysis of the histograms gives a good indication that the transformed data is approximately normally distributed. The histograms have a skewness of approximately zero.

The D'Agostino and Pearson normality test confirms the visual interpretation. The data is normally distributed after Box-Cox with a confidence interval of 95%. The p-values (0.164, 0.982 and 0.354) of all three samples are higher than the significance

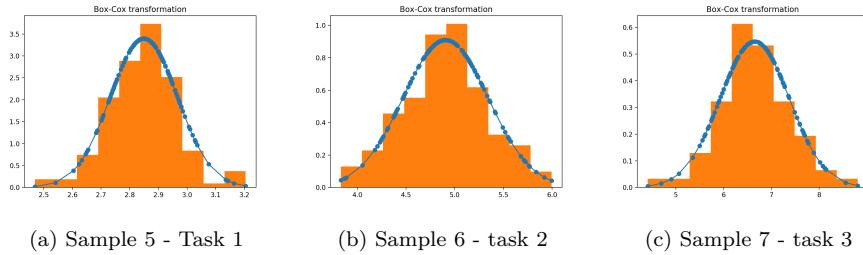


Figure 4.8: Histogram with normal distribution fit after Box-Cox transformation, sample with total time per task

level (0.05). Sample 5, 6 and 7 are normally distributed after the transformation, and the assumptions of normality are met.

4.3.2.4 All participants divided by task, correct element variable

In this section sample 8, 9 and 10 will be normal distribution examined. These samples will be used to test whether there is a significant difference between the three tasks when looking at the number of correctly chosen elements variable.

A visual analysis of the three histograms in figure 4.9a, 4.9b and 4.9c show a negative skewness, just like the histograms in section 4.3.2.2. This give an indication that the three samples are not normally distributed.

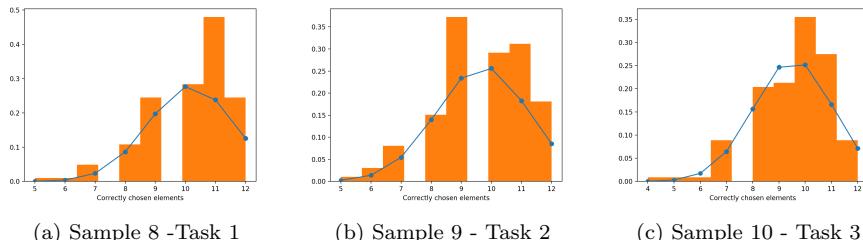


Figure 4.9: Histogram with normal distribution fit showing samples with number of correct elements per task

D'Agostino and Pearson normality test confirms our visual analysis of the histograms in two of three samples. Sample 9 passes the normality test, even though the p-value (0.099) is close to the significance level (0.05). Sample 8 and sample 10 do not pass the normality test with a confidence interval of 95%. Both samples obtained a p-value (0.00022 and 0.0047) smaller than the significant level. The null hypothesis is rejected for sample 8 and 10, and the alternative hypothesis is accepted. The null hypothesis is accepted for sample 9. The significance level is 5%.

4. RESULT

A Box-Cox power transformation is applied to all three samples. The transformation changes the data, and to correctly compare the results, sample 9 also has to be transformed even though it's original data is normally distributed. The transformed samples are shown in histogram 4.10a, 4.10b and 4.10c. All three are negatively skewed, sample 9 and 10 less than sample 8. A visual conclusion is difficult in this case.

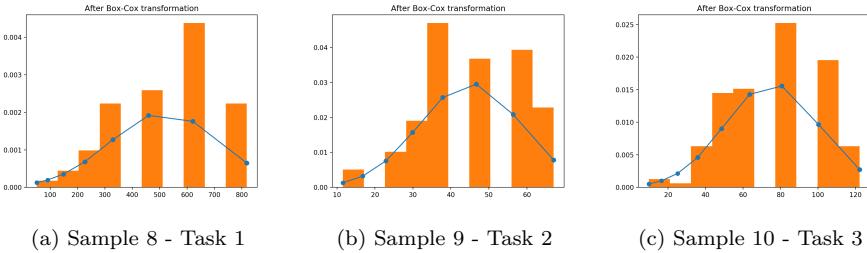


Figure 4.10: Histogram with normal distribution fit after Box-Cox

D'Agostino and Pearson normality test accepts the null hypothesis on sample 9 and 10 and rejects it on sample 8. Sample 9 and 10 has p-values (0.0752 and 0.2104) higher than the significance level, while sample 8's p-value (0.2104) is significantly lower. When using these three samples in hypothesis tests, a non-parametric method should be used. This is because sample 8 is not normally distributed.

4.3.2.5 Experienced participants divided by task, total time variable

In this section, sample 11, 12, and 13, shown in table 4.5, is normal distribution tested. The three samples will be used to test whether there is a significant difference between the three tasks total time results when considering only experienced participants.

A visual interpretation of the histograms in 4.11 show that all three samples are positively skewed (figure 4.2). Skew gives a fairly strong evidence that the samples are not normally distributed.

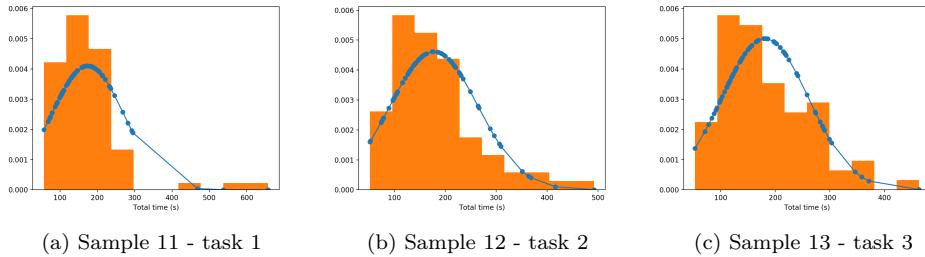


Figure 4.11: Histograms with normal distribution fit with samples containing total time to complete each task

D'Agostino and Pearson normality test confirms our visual interpretation of the three histograms. All three p-values ($1.229 * 10^{-14}$, $2.678 * 10^{-5}$ and 0.000884) are lower than the significance level (5%). Sample 11, 12 and 13 do not pass the normality test with a confidence interval of 95%. The null hypothesis is rejected.

A Box-Cox power transformation is applied to all three samples since the null hypothesis was rejected. Histograms with normal distribution fit containing the transformed data is shown in figure 4.12. Visually, the histograms look normally distributed with minimal skewness.

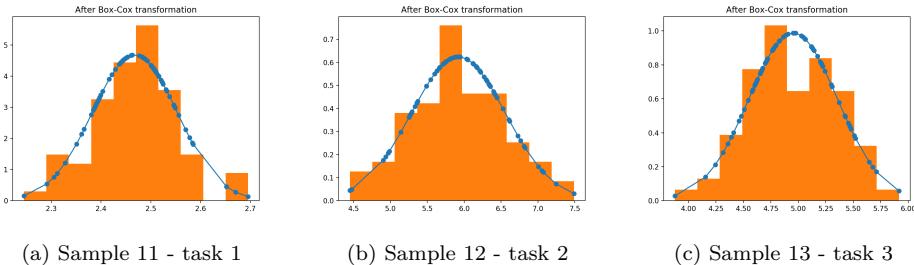


Figure 4.12: Histograms with normal distribution fit containing Box-Cox transformed data

The Box-Cox transformed data is tested with D'Agostino and Pearson normality test. All three samples obtained p-values (0.694, 0.955 and 0.887) larger than the significance level (0.05). Within a confidence interval of 95%, the test concludes that sample 11, 12 and 13 is normally distributed. Sample 11, 12 and 13 can be used in methods assuming normally distributed samples.

4. RESULT

4.3.2.6 Experienced participants divided by task, correct elements variable

This section will test if sample 14, 15 and 16, shown in table 4.6, follows the normal distribution. The three samples will be used to test whether there is a significant difference in the number of correct elements between the three tasks when comparing only experienced participants.

A visual interpretation of the histograms in 4.13 show that all three samples are slightly negatively skewed (figure 4.2). Sample 15 (4.13a) and sample 16 (4.13b) has less skew than sample 14 (4.13c).

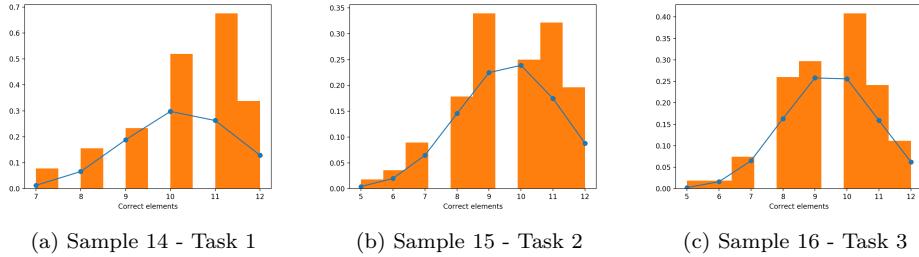


Figure 4.13: Histogram with normal distribution fit showing samples with number of correct elements results for experienced participants

D'Agostino and Pearson normality test confirms our visual interpretation of the three histograms. All three p-values (0.0588, 0.2067 and 0.2975) are higher than the significant level (0.05), notice that sample 14 has a lower p-value than the two other samples. This sample is not as significant as the two other samples. Sample 14, 15 and 16 pass the normality test with a confidence interval of 95%. The null hypothesis is accepted. Sample 14, 15 and 16 can be used in tests that assume normally distributed samples.

4.3.2.7 Inexperienced participants divided by task, total time variable

This section will test if sample 17, 18 and 19, shown in table 4.7, follows the normal distribution. These samples will be used to test whether there is a significant difference in total time between the three tasks when looking at only inexperienced participants.

A visual analysis of the histograms 4.14a, 4.14b and 4.14c show a positive skew. The skew is less than the histograms in figure 4.11, but is most likely too large for the samples to be normally distributed.

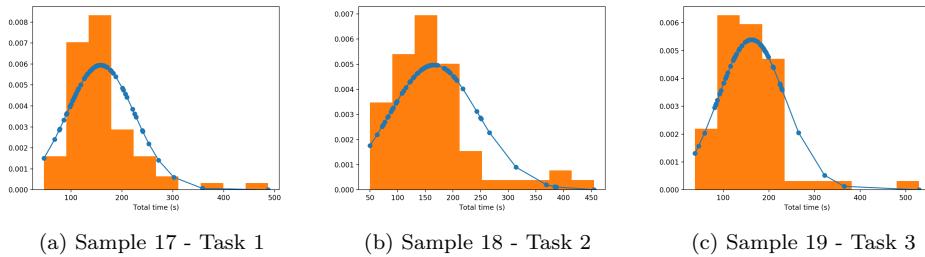


Figure 4.14: Histogram with normal distribution fit

The D'Agostino and Pearson normality test agrees with the visual analysis. The three obtained p-values ($1.586 * 10^{-11}$, $1.773 * 10^{-6}$ and $2.312 * 10^{-11}$) are all significantly lower than the significance level of 0.05. Sample 17, 18 and 19 is not normally distributed with a confidence interval of 95%. The null hypothesis is rejected.

The samples are Box-Cox power transformed. The histograms after transformation (4.15a, 4.15b and 4.15c) are visually evaluated to be normal distributed.

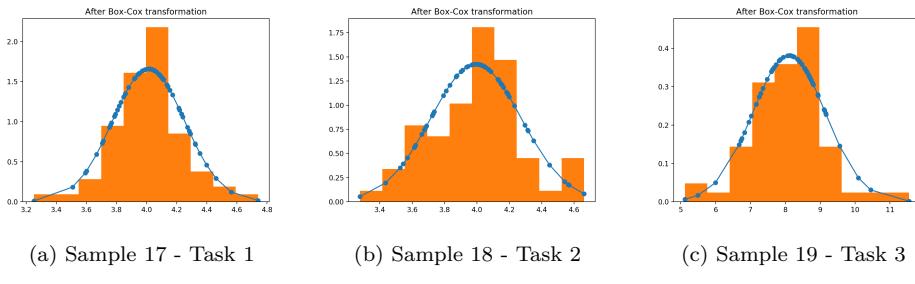


Figure 4.15: Histogram with normal distribution fit after Box-Cox

A new D'Agostino and Pearson normality test on the transformed data confirms that all three samples are normally distributed with a significance level of 5%. The obtained p-values (0.139, 0.909 and 0.067) is smaller than 0.05, and the null hypothesis is accepted. Sample 17, 18, and 19 are normally distributed after Box-Cox and can be used in parametric methods

4. RESULT

4.3.2.8 Inexperienced participants divided by task, correct elements variable

Sample 20, 21 and 22, shown in table 4.8, contains the number of correctly chosen elements in each of the three tasks from only inexperienced participants. The samples will be used to test if inexperienced participants do better in one of the tasks.

A visual interpretation of histogram 4.16a, 4.16b and 4.16c show a negative skew, similar to the histograms containing results from only experienced participants (??).

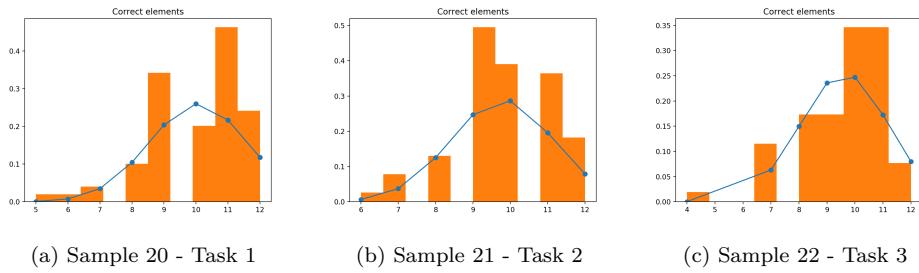


Figure 4.16: Histogram with normal distribution fit

The D'Agostino and Pearson normality test rejects the null hypothesis on sample 20 and 22 and accepts the null hypothesis on sample 21. Obtained p-values (0.007 and 0.004) are lower than 0.05 for sample 20 and 22 and sample 21 (0.523) higher than 0.05. The test concludes that sample 21 is normally distributed, and sample 20 and 22 is not with a significant level of 5%.

A Box-Cox power transformation is applied to all three samples. The transformation changes the data, and to correctly compare the data, sample 21 also has to be transformed even though the original data was followed the normal distribution. The transformed samples are shown in histogram 4.17a, 4.17b and 4.17c. All three histograms are less skewed than the original histograms (4.16).

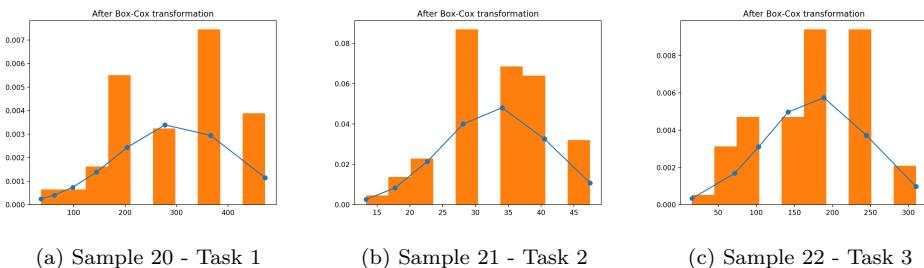


Figure 4.17: Histogram with normal distribution fit after Box-Cox transformation

The D'Agostino and Pearson normality method is executed on the transformed samples. The three obtained p-values (0.061, 0.714 and 0.311) is lower than 0.05. The null hypothesis is accepted. Sample 20, 21 and 22 are normally distributed with a significant level of 5% and can be used in parametric methods.

4.3.2.9 Normality test summary

Table 4.9: Summary of normality tests done in section 4.3.2

	Sample	Normally distributed	Normally distributed after Box-Cox
<i>Total time</i>			
Experienced	1	No	Yes
Inexperienced	2	No	Yes
<i>Correct elements</i>			
Experienced	3	No	No
Inexperienced	4	No	No
<i>Total time</i>			
Task 1	5	No	Yes
Task 2	6	No	Yes
Task 3	7	No	Yes
<i>Correct elements</i>			
Task 1	8	No	No
Task 2	9	Yes	Yes
Task 3	10	No	Yes
<i>Total time, experienced participants</i>			
Task 1	11	No	Yes
Task 2	12	No	Yes
Task 3	13	No	Yes
<i>Correct elements, experienced participants</i>			
Task 1	14	Yes	<i>not tested</i>
Task 2	15	Yes	<i>not tested</i>
Task 3	16	Yes	<i>not tested</i>
<i>Total time, inexperienced participants</i>			
Task 1	17	No	Yes
Task 2	18	No	Yes
Task 3	19	No	Yes
<i>Correct elements, inexperienced participants</i>			
Task 1	20	No	Yes
Task 2	21	Yes	Yes

4. RESULT

Task 3	22	No	Yes
--------	----	----	-----

4.3.3 Levene's test of Equality of Variance

As mentioned in section 4.2.2.1, 4.2.2.2 and 4.2.2.4, the two sample t-test, one-way ANOVA and Mann-Whitney U test assumes that the samples come from populations with equal variances. This assumption will be examined with Levene's test. The hypothesis tested is:

H_0 : Input samples are from populations with equal variances

H_A : Input samples are from populations that do not have equal variances

The null hypothesis is accepted if the obtained p-value is higher than the significance level. Table 4.10 contains the summary of the Levene's test performed on all sample pairs. All sample pairs accepted the null hypothesis except sample 1 and 2, who obtained a p-value lower than the significance level. The test used a significance level of 5% on all the tests.

Table 4.10: Summary of Levene's tests

Sample	Obtained p-value	Samples are from populations with equal variances
<i>Total time, all</i> 1 and 2	0.030	No
<i>Correct elements, all</i> 3 and 4	0.823	Yes
<i>Total time, divided by task</i> 5, 6, and 7	0.636	Yes
<i>Correct elements, divided by task</i> 8, 9, and 10	0.805	Yes
<i>Total time, experienced participants divided by task</i> 11, 12, and 13	0.972	Yes
<i>Correct elements, experienced participants divided by task</i> 14, 15, and 16	0.724	Yes
<i>Total time, inexperienced participants divided by task</i> 17, 18, and 19	0.499	Yes
<i>Correct elements, inexperienced participants divided by task</i>		

20, 21, and 22	0.626	Yes
----------------	-------	-----

4.3.4 Hypothesis testing

This section will test all the hypothesis. The test order will be the same as the tables in section 4.3.1. Which statistic test that is used is determined by the results from the normality test (section 4.3.2) and equal variance test (section 4.3.3).

4.3.4.1 Test differences in total time between experienced and inexperienced participants

This section will test if there are any difference in total time spent on the tasks between experienced and inexperienced participants. This test is covered by sample 1 and sample 2 from section ???. Sample 1 is experienced participants and sample 2 inexperienced participants. Both samples was normally distributed after a Box-Cox transformation (4.3.2.1). A two-sample t-test will be used to answer this hypothesis since the normality assumption is valid. The hypothesis tested in this section is:

$$H_0: \text{Equal task time between experienced and inexperienced participants}$$

$$H_A: \text{Unequal task time between experienced and inexperienced participants}$$

If \bar{x}_1 equals the mean time for experienced-, and \bar{x}_2 the mean time for inexperienced participants, the hypothesis can be written as:

$$\begin{aligned} H_0: \bar{x}_1 &= \bar{x}_2 \\ H_A: \bar{x}_1 &\neq \bar{x}_2 \end{aligned}$$

Since we cannot assume equal variances in the two samples (Table 4.10), this test will use the Welch's t-test for unequal variances [(Walpole et al., 2012), p. 345]. Equation 4.1 is still valid. The obtained values from the test are shown in the gray box. The obtained T-statistic is smaller than the critical value. The t-test then conclude that there is a significant difference between the means of the two population samples with a confidence interval of 95%.

4. RESULT

Two sample, two-way t-test, sample 1 and 2
Significance level: 5%

T – statistic: -60.442
Degree of freedom (v): 447
Significance level (α): 0.05
Critical value: 1.960

Using equation 4.1, the absolute value of the *T – statistic* is larger than the critical value ($|60.442| > 1.960$) and the null hypothesis is rejected and H_A accepted.

Test if experienced or inexperienced participants finish the task fastest

Because there was a statistical significant difference between time spent on each task between the participants, this section will also test who finished the task fastest. The second hypothesis tested in this section is:

$$H_0: \text{Equal task time between all participants}$$
$$H_A: \text{Experience participants finish the task faster}$$

With sample 1 equals experienced participants and sample 2 equals inexperienced participants we get the hypothesis:

$$H_0: \bar{x}_1 = \bar{x}_2$$
$$H_A: \bar{x}_1 < \bar{x}_2$$

This test gives the same T-statistics as the previous test, but the critical value is changed since this test used in the second hypothesis is a two sample, one-way t-test. Since we cannot assume equal variances in the two samples a Welch's test is performed. Obtained *T – statistic* is still smaller than the critical value ($-64.654 < 1.645$). Our test is to check if sample 1's mean is significantly larger than sample 2's mean, then we use the comparison test written in equation 4.2, section 4.2.2.1. Our *T – statistic* is not larger than the critical value, and we need to accept the null hypothesis. There is no evidence that experienced participants use less time on the tasks than the inexperienced.

Two sample, one-way t-test, sample 1 and 2
Significance level: 5%

T – statistic: -60.442
Degree of freedom (v): 447
Significance level (α): 0.05
Critical value: 1.645

T-statistic is smaller than the critical value ($-60.442 < 1.645$) and the null hypothesis is accepted.

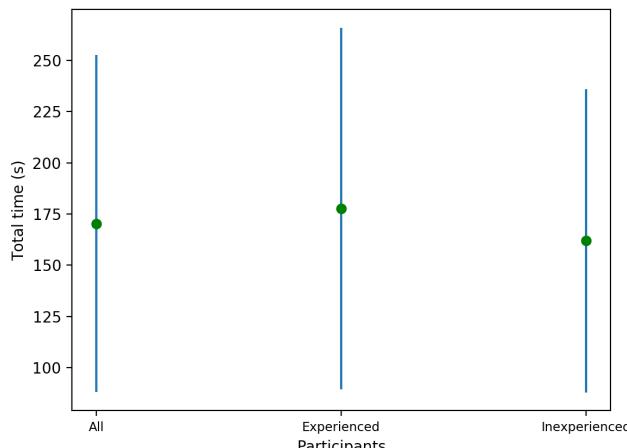


Figure 4.18: Sample 1 and 2 - mean (green dot) and standard deviation (blue line)

Since we know that there is a statistical significant difference between the two sample means, the author concludes that the inexperienced participants finished the task faster than the experienced participants. The time difference can also be seen in plot 4.18. Inexperienced participants finished the tasks in average 16 seconds faster than experienced participants.

4.3.4.2 Test if there is a difference between experienced and inexperienced participants in total correct elements

This section will test if there is a difference between experienced- and inexperienced participants when looking at the number of correctly chosen elements. Sample 3 and 4 is the correct samples to use in this test. Both samples are not normally distributed,

4. RESULT

and we need to use a non-parametric method. The Mann-Whitey U test is the preferred test to use on these samples. As mentioned in section 4.2.2.3, the Mann-Whitey U test is preferred when the samples have ties (identical observations). From histogram 4.6a and 4.6b we see that the samples are identical in some cases. Mann-Whitey U test should, therefore, be used to compare the population medians. The hypothesis to be tested is:

$$\begin{aligned} H_0: \text{median}_3 &= \text{median}_4 \\ H_A: \text{median}_3 &\neq \text{median}_4 \end{aligned}$$

Using equation 4.5 in section 4.2.2.4 and the obtained $U - \text{statistic}$, we conclude that there is not enough evidence to reject the null hypothesis with a confidence interval of 95%. The $U - \text{statistic}$ is larger than the critical value, and the null hypothesis is accepted.

Two sample t-test, sample 3 and 4
Significance level: 5%

$U - \text{statistic}$: 17012
Significance level (α): 0.05
Sample size, n1: 229
Sample size, n2: 200
Critical-value: 127

$U - \text{statistic}$ is larger than the critical value ($17012 > 127$)
and the null hypothesis is accepted.

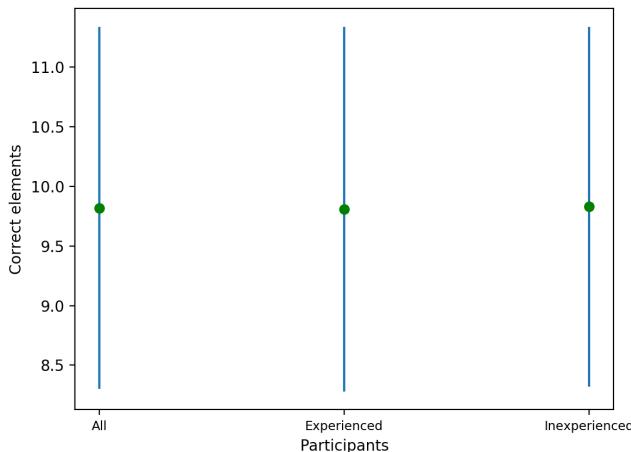


Figure 4.19: Sample 3 and 4 - mean (green dot) and standard deviation (blue line)

Results from this section show that there is not enough evidence to conclude that there is any difference between experienced- and inexperienced participants when looking at the number of correctly chosen elements per task. The author concludes that experienced- and inexperienced participants did equally well on the task. This can also be seen visually in figure 4.19. Mean values show in the figure is approximately equal between the participants.

4.3.4.3 Test if total time differs between the three tasks

This section will test if time spent varies between each of the three tasks. Sample 5, 6, and 7 is used in this test. The one-way *ANOVA* test will be used in this section. This test is used when more than two samples are compared (section 4.2.2.2). The assumption that sample 5, 6 and 7 come from populations with equal variances are met 4.10), the samples are also normally distributed after a Box-Cox transformation (4.3.2.3). The hypothesis tested here is:

$$H_0: \bar{x}_5 = \bar{x}_6 = \bar{x}_7$$

H_A : Total time differs between at least two of the tasks

Using equation 4.4 in section 4.2.2.2 and results obtained from the calculations, the one-way *ANOVA* test rejects the null hypothesis. The obtained *f – value* is lower than the critical value. With a confidence interval of 95% the author claim that there is a difference between the mean value of the three tasks.

4. RESULT

One-way ANOVA, sample 5, 6 and 7
Significance level: 5%

f – value: 2123.308
Significance level (α): 0.05
 $v_1 = 2, v_2 = 426$
Critical-value: 3.00

f – value is significantly higher than the critical value (2123.308 > 3.00) and the null hypothesis is rejected, H_A is accepted

Since the null hypothesis was rejected *Tukey's method*, is used to make comparisons between task one, two, and three. This test did not find any statistically significant difference between the three tasks. Visual evaluation of figure 4.20 show that task 1 was completed slightly faster than the two other tasks.

Tukey's test, sample 5, 6 and 7
Significance level: 5%

Task 1 - Sample 5, Task 2 - Sample 6 and Task 3 - Sample 7

Task 1 and Task 2 do not differ significantly
Task 1 and Task 3 do not differ significantly
Task 2 and Task 3 do not differ significantly

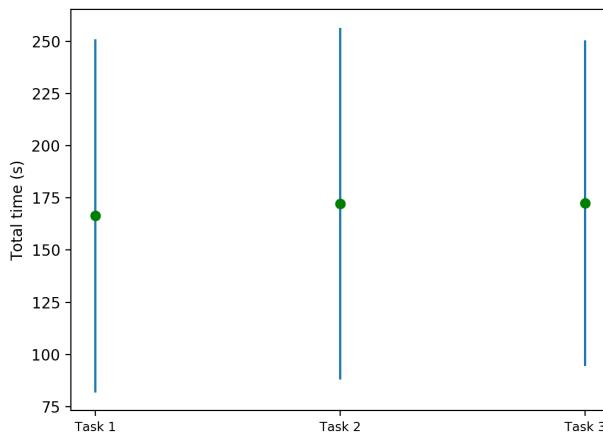


Figure 4.20: Sample 5, 6 and 7 - mean (green dot) and standard deviation (blue line)

Based on the results from this section the author concludes that there is a statistically significant difference in total time between at least two of the tasks, but the difference is not enough for *Tukey's test* to find a significant difference.

4.3.4.4 Test if the number of correct elements differs between the three tasks

This section will test if there is a difference in the number of correctly chosen elements between the three tasks. The test will use sample 8, 9 and 10. Since sample 8 are not normally distributed (4.3.2.4) a non-parametric test should be applied. The Kruskal-Wallis test is the non-parametric equivalent to one-way ANOVA (4.2.2.5). It is used to test equality of medians when the samples are not normally distributed. The hypothesis tested is:

$$H_0: \text{median}_8 = \text{median}_9 = \text{median}_{10}$$

H_A : Number of correctly chosen elements differ between at least two of the tasks

Using equation 4.6 in section 4.2.2.5, the Kruskal-Wallis test rejects the null hypothesis. The obtained H – value is smaller than the critical value. The p-value is approximately zero, and this gives a good indication that the result is significant. With a confidence interval of 95% the author claim that there is a difference between the median value of the three tasks.

4. RESULT

Kruskal-Wallis test, sample 8, 9 and 10
Significance level: 5%

P – value: $3.967 * 10^{-72}$

H – value: 328.816

Significance level (α): 0.05

$v = 2$

Critical-value: 5.991

H – value is significantly higher than the critical value
($328.816 > 5.991$) and the null hypothesis is rejected, H_A is accepted

Like in one-way ANOVA, a *post hoc* test should be used to make paired comparisons to determine which groups differ. The *post hoc* test applied is Tukey's test. Results from Tukey's test resulted in a significant difference in the number of correctly chosen elements between task 1 and task 2, and task 1 and task 3. Figure 4.21 show that task 1 has a higher mean value than the other two tasks. Task 1 also has a smaller standard deviation than the other tasks.

Tukey's test, sample 8, 9 and 10
Significance level: 5%

Task 1 - Sample 8, Task 2 - Sample 9 and Task 3 - Sample 10

Task 1 and Task 2 differs significantly

Task 1 and Task 3 differs significantly

Task 2 and Task 3 do not differ significantly

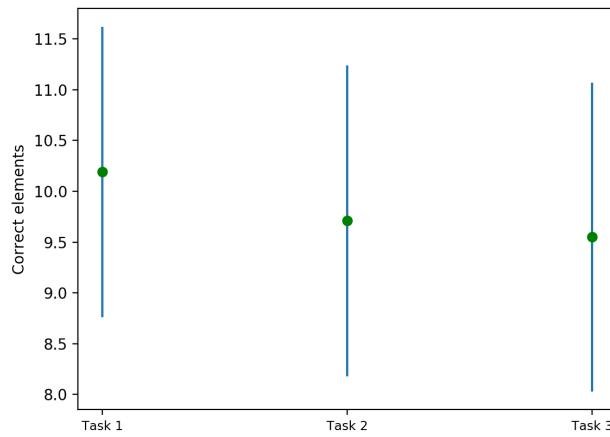


Figure 4.21: Sample 8, 9 and 10 - mean (green dot) and standard deviation (blue line)

The results in this section find statistically significant evidence that the participants had a higher number of correct elements in Task 1 than in the two other tasks.

4.3.4.5 Test differenced in results from experienced participants

This section will answer two hypothesis about experienced participant's results divided by task. The first hypothesis will test if time spent on each of the three tasks differ and the second hypothesis will test if the number of correct elements in each of the three tasks differs when the samples only include results from experienced participants. Sample 11, 12 and 13 will be used to answer the first hypothesis. These three samples are normally distributed after a Box-Cox transformation (4.9) and also come from populations with equal variances (??). Sample 14, 15 and 16 will be used on the second hypothesis. All three samples are normally distributed (4.9) and also come from populations with equal variances (??). The one-way *ANOVA* method will be used to test both hypotheses.

The first hypothesis is:

$$H_0: \bar{x}_{11} = \bar{x}_{12} = \bar{x}_{13}$$

$$H_A: \text{Total time is different between at least two of the tasks}$$

Using equation 4.4 from section 4.2.2.2, the one-way *ANOVA* test rejects the null hypothesis. The obtained *f-value* (1216.919) is higher than the critical value (3.00). The calculated p-value is also approximately zero, and this gives a good indication

4. RESULT

that the result is significant. With a confidence interval of 95% the author claim that there is a time difference between the three tasks.

When the null hypothesis is rejected, a *post hoc* test is used to compare each task with each other. Tukey's *post hoc* test did not find any significant difference between the three tasks with a significant level of 5%. Figure 4.22 show an approximately similar mean value in all three tasks.

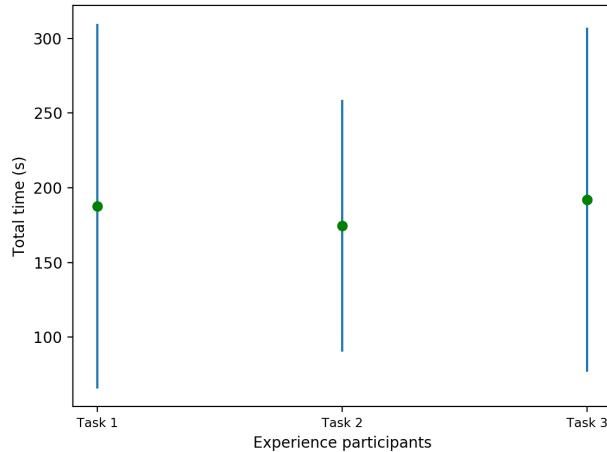


Figure 4.22: Sample 11, 12 and 13 - mean (green dot) and standard deviation (blue line)

The second hypothesis is:

$$H_0: \bar{x}_11 = \bar{x}_12 = \bar{x}_13$$

H_A : Number of correct elements in each task differs between at least two of the tasks

Using equation 4.4, the one-way ANOVA test rejects the null hypothesis. The obtained *f – value* (8.210) is higher than the critical value (3.00). The p-value is also approximately zero and this gives a good indication that the result is significant. With a confidence interval of 95% the author claim that there is a difference between the mean value of at least two of the tasks.

Since the null hypothesis was rejected, a *post hoc* test should be used to make paired comparisons to determine which groups differ. Tukey's *post-hoc* test resulted in a significant difference in the number of correctly chosen elements between task 1 and task 2, and task 1 and task 3. Figure 4.23 show that task 1 has a higher mean value than the two other tasks. Task 1 also has a smaller standard deviation.

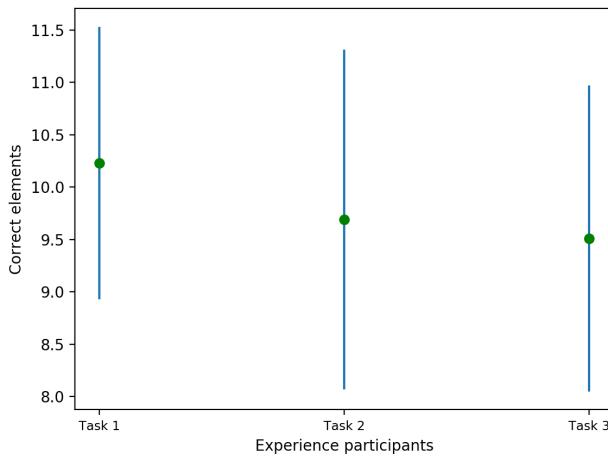


Figure 4.23: Sample 14, 15 and 16 - mean (green dot) and standard deviation (blue line)

With the results found in this section, the author concludes that there is a statistically significant difference in the number of correctly chosen elements between the three tasks. The experienced participants got the best result on task 1. Time spent on each task also differ between at least two of the tasks, but the difference is not significant enough so that Tukey's test can determine a difference. Figure 4.22 show that task 2 has the lowest mean time value of the three tasks.

4.3.4.6 Test differenced in results from inexperienced participants

This section will answer the same hypothesis as the previous section, only with results from inexperienced participants. The first hypothesis will test if time spent on each task differ and the second hypothesis will test if the number of correct elements in each task differs. Sample 17, 18 and 19 will be applied in the first hypothesis test. These samples are normally distributed (4.3.2.7) and come from populations with equal variances (??). Sample 20, 21 and 22 will be used on the second hypothesis. All three samples are normally distributed (4.3.2.8) and also come from populations with equal variances (??). The one-way ANOVA will be used to test both hypotheses.

The first hypothesis is:

$$H_0: \bar{x}_{17} = \bar{x}_{18} = \bar{x}_{19}$$

$$H_A: \text{Total time differ between at least two of the tasks}$$

The one-way ANOVA test rejects the null hypothesis and accepts the alternative hypothesis (H_A) with a significant level of 5%. The obtained f-value from the test is

4. RESULT

higher than the critical value ($905.34 > 3.00$). Since the alternative hypothesis was accepted, Tukey's *post hoc* test is used to make compared comparisons between task 1, task 2 and task 3. The test does not find a significant difference when comparing each of the three tasks with a significant level of 5%. Figure 4.24 show that inexperienced participants spent more time on task 2 than the other tasks, but the difference is not significant according to Tukey's test.

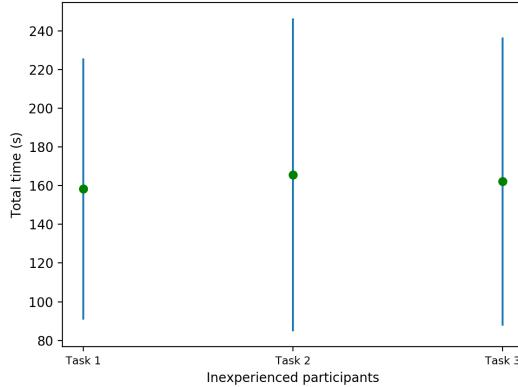


Figure 4.24: Mean (green dot) and standard deviation (blue line) for sample 17, 18 and 19

The seconds hypothesis is:

$$H_0: \bar{x}_{20} = \bar{x}_{21} = \bar{x}_{22}$$

H_A : Number of correct elements differ between at least two of the tasks

The one-way ANOVA test rejects the null hypothesis (H_0) and accepts the alternative hypothesis (H_A) with a significance level of 5%. The obtained f-value from the test is higher than the critical value ($189.05 > 3.00$). Since the null hypothesis was rejected, Tukey's *post hoc* test will be used to make comparisons between the three tasks. This test cannot find a significant difference when comparing the tasks with a significant level of 5%. Looking at figure 4.25, task 1 has a higher mean time value than the two other tasks, task 2 also has a higher mean than task 3. Even though there are differences in the number of correct elements between the tasks, it is not significant according to Tukey's test.

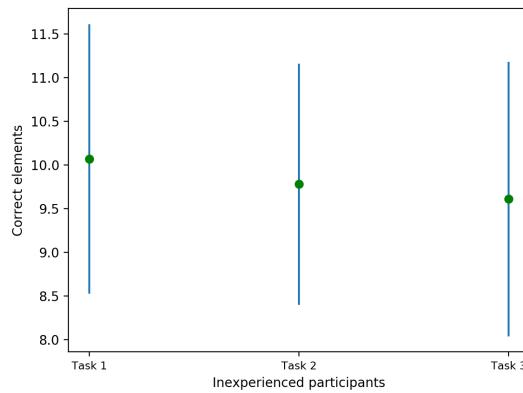


Figure 4.25: Mean (green dot) and standard deviation (blue line) for sample 20, 21 and 22

With the results found in this section, the author concludes that there is a statistically significant difference in both total time spent on each task and the number of correctly chosen elements between at least two of the tasks. The differences are not significant enough so that Tukey's test can determine which task differs. Figure 4.24 and 4.25 show that task 1 has the lowest mean time value and the highest mean correct value.

4.3.4.7 Hypothesis test summary

Table 4.11: Summary of hypothesis tests done in section 4.3.4

Hypothesis (Dependent variable, Independent variable)	Participants Sample number	Hypothesis is accepted
Total time, Experienced and Inexperienced There are a difference between experienced and inexperienced participants	All 1 and 2	Yes
Experienced participants finish the tasks faster than inexperienced	1 and 2	No
Inexperienced participants finish the tasks faster than experienced	1 and 2	Yes
Correct elements, Experienced and Inexperienced There are a difference between experienced and inexperienced participants	All 3 and 4	No
Total time, Task 1, Task 2 and Task 3 Total time is different between at least two of the tasks Task 1 significantly differs from Task 2	All 5, 6 and 7 5, 6 and 7	Yes No

4. RESULT

Task 1 significantly differs from Task 3	5, 6 and 7	No
Task 2 significantly differs from Task 3	5, 6 and 7	No
Correct elements, Task 1, Task 2, Task 3	<i>All</i>	
The number of correctly chosen elements is different between at least two of the tasks	8, 9 and 10	Yes
Task 1 significantly differs from Task 2	8, 9 and 10	Yes
Task 1 significantly differs from Task 3	8, 9 and 10	Yes
Task 2 significantly differs from Task 3	8, 9 and 10	No
Total time, Task 1, Task 2, Task 3	<i>Experienced</i>	
Total time is different between at least two of the tasks	11, 12 and 13	Yes
Task 1 significantly differs from Task 2	11, 12 and 13	No
Task 1 significantly differs from Task 3	11, 12 and 13	No
Task 2 significantly differs from Task 3	11, 12 and 13	No
Correct elements, Task 1, Task 2, Task 3	<i>Experienced</i>	
Number of correct elements differs between at least two of the tasks	14, 15 and 16	Yes
Task 1 significantly differs from Task 2	14, 15 and 16	Yes
Task 1 significantly differs from Task 3	14, 15 and 16	Yes
Task 2 significantly differs from Task 3	14, 15 and 16	No
Total time, Task 1, Task 2, Task 3	<i>Inexperienced</i>	
Total time is different between at least two of the tasks	17, 18 and 19	Yes
Task 1 significantly differs from Task 2	17, 18 and 19	No
Task 1 significantly differs from Task 3	17, 18 and 19	No
Task 2 significantly differs from Task 3	17, 18 and 19	No
Correct elements, Task 1, Task 2, Task 3	<i>Inexperienced</i>	
Number of correct elements differs between at least two of the tasks	20, 21 and 22	Yes
Task 1 significantly differs from Task 2	20, 21 and 22	No
Task 1 significantly differs from Task 3	20, 21 and 22	No
Task 2 significantly differs from Task 3	20, 21 and 22	No

5 | Discussion

Task 1 had an average difficulty of 2.11, task 2 of 2.14 and task 3 of 2.5 (task 4 1.79). Experienced participants gave a lower difficulty score than inexperienced on all three tasks (in average a 0.20 lower difficulty score).

6 | Proposed sections

6.1 Future work

Create a survey to test how accurate both experienced and inexperienced participants digitize buildings from aerial images. Can use FKB as the correct polygon and compare it with the drawn polygon from participants.

Do a study with reward. Compare reward and not reward geo tasks. Do they solve the tasks better with reward? "A reward can be provided for merely participating in the task. The reward can also be provided as a prize for submitting the best solution or one of the best solutions. Thus, the reward can provide an incentive for members of the community to complete the task as well as to ensure the quality of the submissions."

The future in micro-tasking "belongs to hybrid methodologies that combine human computation with advanced computing" (Meier, 2013b).

When aiming towards wider adoption of crowdsourcing one have to be aware of the challenges of using it. It is important to remember that all tasks do not fit into the micro-tasking crowd worker model. Very complex tasks that can't be partitioned are not suitable for solving through micro-tasks.

Advanced computing techniques such as Artificial Intelligence and Machine Learning is needed to build approaches that combine the power of people with the speed and scalability of automated algorithms (Meier, 2013b).

6.2 Usage potential

Systems are exploiting the people's physical presence in an environment more, they are more location dependent. This can be particularly important when seeking to improve geospatial data quality [(Meier, 2013b), p. 323]. "For instance, UrbanMatch (Celino et al. 2012a) is a mobile location based game that uses player's familiarity with a city to link photos with points of interest in the city. Players are shown points of interest and known images from a trusted source (e.g. OpenStreetMap) and asked if photos from an untrusted source (e.g. Flickr) might also relate to the point of interest".

(Meier, 2013b): "As the previous sections show there is a lot of potential for AR systems to use HC to provide content, and to support processing in other ways. However there has been little research to date combining AR and HC systems. In this section we review the first research efforts in this area."

(Meier, 2013b) "Lastly, there is huge untapped potential in leveraging the "cognitive surplus" available in massively multiplayer online games to process humanitarian mi-

6. PROPOSED SECTIONS

crotasks during disasters. The online game “League of Legends,” for example, has 32 million players every month and three million on any given day. Over 1 billion hours are spent playing League of Legends every month. Riot Games, the company behind League of Legends is even paying salaries to select League of Legend players. Now imagine if users of the game were given the option of completing microtasks in order to acquire additional virtual currency, which can buy better weapons, armor, etc. Imagine further if users were required to complete a microtask in order to pass to the next level of the game. Hundreds of millions of humanitarian microtasks could be embedded in massively multiplayer online games and instantaneously completed. Maybe the day will come when kids whose parents tell them to get off their computer game and do their homework will turn around and say: “Not now, Dad! I’m microtasking crisis information to help save lives in Haiti!” ”

Machines are bad at tackling things they have never seen before. They need to learn from large amounts of passed data. Humans don’t need this. Humans can solve tasks we have never seen before. Tackling new/novel situations are humans much better than machines. Business strategies, marketing holes, this are tasks only humans can do.

Data Categorization, organize your data, no matter what the data is. Micro-tasking platforms can turn all the big data into rich data that is organized, streamlined, and useful. Micro-tasking let’s you organize your original data which again can be used to train machine learning models. According to CrowdFlower is human-curated training sets the best traning datasets to use.

Appendices

A | Tets

Fbox

Some text esfljsf
lksj lksdjflsk slk

Some text
kduhaszkdh aszkd-
jhs zkjdfh skdj
skd

dwkjdkwjd dh wkjdhw kjdh wkjhd qwkjhd kwd qw .

text

dwkjdkwjd dh wkjdh wkjhd qwkjhd kwd qw .

Bibliography

- Affairs, A. S. f. P. (2013). System Usability Scale (SUS).
- Bangor, A., Kortum, P., and Miller, J. (2009). Determining What Individual SUS Scores Mean: Adding an Adjective Rating Scale. *Journal of Usability Studies*, 4(3):114–123.
- Barrouillet, P., Bernardin, S., Portrat, S., Vergauwe, E., and Camos, V. (2007). Time and Cognitive Load in Working Memory.
- Ben, S. and Plaisant, C. (2009). *Designing the User Interface*. Pearson, fifth edition.
- Bernstein, M. S., Little, G., Miller, R. C., Hartmann, B., Ackerman, M. S., Karger, D. R., Crowell, D., and Panovich, K. (2015). Soylent: A Word Processor with a Crowd Inside. *COMMUNICATIONS OF THE ACM*, 58(8):85–94.
- Biewald, L. (2015). Why human-in-the-loop computing is the future of machine learning | Computerworld. Date accessed: 2017-05-14 URL: <http://www.computerworld.com/article/3004013/robotics/why-human-in-the-loop-computing-is-the-future-of-machine-learning.html>.
- Brooke, J. (1996). *SUS-A quick and dirty usability scale. "Usability Evaluation In Industry"*. Taylor & Francis.
- Deng, X., Joshi, K. D., and Galliers, R. D. (2016). THE DUALITY OF EMPOWERMENT AND MARGINALIZATION IN MICROTASK CROWDSOURCING: GIVING VOICE TO THE LESS POWERFUL THROUGH VALUE SENSITIVE DESIGN 1. 40(2):279–300.
- Difallah, D. E., Catasta, M., Demartini, G., Ipeirotis, P. G., and Cudré-Mauroux, P. (2015). The Dynamics of Micro-Task Crowdsourcing. *Proceedings of the 24th International Conference on World Wide Web - WWW '15*, pages 238–247.
- Difallah, D. E., Demartini, G., and Cudré-Mauroux, P. (2016). Scheduling Human Intelligence Tasks in Multi-Tenant Crowd-Powered Systems. *Proceedings of the 25th International Conference on World Wide Web - WWW '16*, pages 855–865.
- EYeka (2015). The State of crowdsourcing 2015 - How the world's biggest brands and companies are opening up to consumer creativity. Technical report.
- Fan, H., Zipf, A., Fu, Q., and Neis, P. (2014). Quality assessment for building footprints data on OpenStreetMap. *International Journal of Geographical Information Science*, (28:4):700–719.
- Franklin, M. J., Kossmann, D., Kraska, T., Ramesh, S., and Xin, R. (2011). CrowdDB: answering queries with crowdsourcing. *Proceedings of the 2011 international conference on Management of data - SIGMOD '11*, pages 61–72.

BIBLIOGRAPHY

- Frost, J. (2015). Choosing Between a Nonparametric Test and a Parametric Test. Date accessed: 2017-04-23 URL: <http://blog.minitab.com/blog/adventures-in-statistics-2/choosing-between-a-nonparametric-test-and-a-parametric-test>.
- Gadiraju, U., Demartini, G., Kawase, R., and Dietze, S. (2015). Human Beyond the Machine: Challenges and Opportunities of Microtask Crowdsourcing. *IEEE Intelligent Systems*, 30(4):81–85.
- Ghasemi, A. and Zahediasl, S. (2012). Normality tests for statistical analysis: a guide for non-statisticians. *International journal of endocrinology and metabolism*, 10(2):486–9.
- Ha, A. (2016). CrowdFlower raises \$10M to combine artificial intelligence with crowdsourced labor | TechCrunch. Date Accessed: 2017-05-08 URL: <https://techcrunch.com/2016/06/07/crowdflower-series-d/>.
- Holzinger, A. (2013). Human–Computer Interaction and Knowledge Discovery (HCI-KDD): What Is the Benefit of Bringing Those Two Fields to Work Together? *Springer Lecture Notes in Computer Science LNCS 8127*, pages 319–328.
- Holzinger, A. (2016). Interactive machine learning for health informatics: when do we need the human-in-the-loop? *Brain Informatics*, 3(2):119–131.
- Holzinger, A., Plass, M., Holzinger, K., Crișan, G. C., Pintea, C.-M., and Palade, V. (2016). Towards interactive Machine Learning (iML): Applying Ant Colony Algorithms to Solve the Traveling Salesman Problem with the Human-in-the-Loop Approach. In *Availability, Reliability, and Security in Information Systems*, pages 81–95. Springer, Cham.
- Howe, J. (2006). The Rise of Crowdsourcing. *Wired Magazine*, 14(06):1–5.
- Huotari, K. and Hamari, J. (2017). A definition for gamification: anchoring gamification in the service marketing literature. *Electronic Markets*, 27(1):21–31.
- Ipeirotis, P. G. and G., P. (2010). Analyzing the Amazon Mechanical Turk marketplace. *XRDS: Crossroads, The ACM Magazine for Students*, 17(2):16–21.
- ISO (1998). ISO 9241-11:1998(en), Ergonomic requirements for office work with visual display terminals (VDTs) — Part 11: Guidance on usability.
- Israel, G. D. (1992). Determining Sample Size 1.
- Kiefer, P., Giannopoulos, I., Duchowski, A., and Raubal, M. (2016). Measuring Cognitive Load for Map Tasks Through Pupil Diameter. pages 323–337. Springer, Cham.
- Kitchin, R. and Tate, N. J. (2000). *Conducting Research into Human Geography*. Prentice Hall.
- Kostas (2016). Using Crowdsourcing and Machine Learning to locate swimming pools in Australia . Tomnod. Date Accessed: 2017-05-04 URL: <http://blog.tomnod.com/crowd-and-machine-combo>.

-
- LaMorte, W. W. (2017). Mann Whitney U Test (Wilcoxon Rank Sum Test). Date Accessed: 2017-05-12 URL: http://sphweb.bumc.bu.edu/otlt/mpb-modules/bs/bs704_nonparametric/BS704_Nonparametric4.html.
- Leppink, J., Paas, F., Van Gog, T., Van Der Vleuten, C. P. M., and Van Merriënboer, J. J. G. (2014). Effects of pairs of problems and examples on task performance and different types of cognitive load. *Learning and Instruction*, 30:32–42.
- Lund Research Ltd (2013a). Hypothesis Testing - Significance levels and rejecting or accepting the null hypothesis.
- Lund Research Ltd (2013b). Mann-Whitney U Test in SPSS Statistics | Setup, Procedure & Interpretation | Laerd Statistics. Date Accessed: 2017-05-12 URL: <https://statistics.laerd.com/spss-tutorials/mann-whitney-u-test-using-spss-statistics.php>.
- Lund Research Ltd (2013c). One-way ANOVA - Its preference to multiple t-tests and the assumptions needed to run this test | Laerd Statistics. Date Accessed: 2017-04-25 URL: <https://statistics.laerd.com/statistical-guides/one-way-anova-statistical-guide-2.php>.
- Mandler, G. (2013). The Limit of Mental Structures, The Journal of General Psychology, pages 243–250.
- MedCalc Software bvba (2017). Skewness and Kurtosis. Date accessed: 2017-04-23 URL: <https://www.medcalc.org/manual/skewnesskurtosis.php>.
- Meier, P. (2013a). Digital Humanitarian Response: Moving from Crowdsourcing to Microtasking | iRevolution. Date Accessed: 2017-05-04 URL: <https://irevolutions.org/2013/01/20/digital-humanitarian-micro-tasking/>.
- Meier, P. (2013b). Handbook of Human Computation. In *Handbook of Human Computation*. Springer New York, New York, NY.
- Meier, P. (2014). Typhoon | iRevolution. Date accessed: 2017-05-07 URL: <https://irevolutions.org/tag/typhoon/>.
- Michelucci, P. and Dickinson, J. L. (2016). The power of crowds. *Science*, 351(6268):32–33.
- Morschheuser, B., Hamari, J., and Koivisto, J. (2016). Gamification in Crowdsourcing: A Review. In *2016 49th Hawaii International Conference on System Sciences (HICSS)*, pages 4375–4384. IEEE.
- Motulsky, H. (2013). Intuitive Biostatistics: A Nonmathematical Guide to Statistical Thinking, 3rd edition. In *Intuitive Biostatistics*, chapter 24.
- Nikki (2016). Finding Swimming Pools in Australia using Deep Learning. Date Accessed: 2017-05-04 URL: <http://blog.tomnod.com/finding-pools-with-deep-learning>.

BIBLIOGRAPHY

- Oppenheimer, D. (2017). Machine Learning with Humans in the Loop - Algorithmia. Date accessed: 2017-05-14 URL: <http://blog.algorithmia.com/machine-learning-with-human-in-the-loop/>.
- Osborne, J. W. (2010). Improving your data transformations: Applying the Box-Cox transformation. *Practical Assessment, Research & Evaluation*, 15(12).
- Palen, L., Soden, R., Anderson, T. J., and Barrenechea, M. (2015). Success & Scale in a Data-Producing Organization. *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems - CHI '15*, pages 4113–4122.
- Pearson, A., Sözcükler, A., düzeltmeli Kolmogorov-Smirnov, L., Pearson ve Jarqua-Bera testleri Derya ÖZTUNA Atilla Halil ELHAN Ersöz TÜCCAR, A., and Öz-tuna, D. (2006). Investigation of Four Different Normality Tests in Terms of Type 1 Error Rate and Power under Different Distributions. *Turk J Med Sci*, 36(3):171–176.
- Quinn, A. J. and Bederson, B. B. (2011). Human computation. *Proceedings of the 2011 annual conference on Human factors in computing systems - CHI '11*, pages 1403–1412.
- Salk, C., Sturn, T., See, L., and Fritz, S. (2016). Local Knowledge and Professional Background Have a Minimal Impact on Volunteer Citizen Science Performance in a Land-Cover Classification Task. *Remote Sensing*, 8(10):774.
- Sarasua, C., Simperl, E., and Noy, N. F. (2012). Crowdsourcing Ontology Alignment with Microtasks. pages 525–541.
- Schade, A. (2015). Pilot Testing: Getting It Right (Before) the First Time. Date accessed: 2017-04-19 URL: <https://www.nngroup.com/articles/pilot-testing/>.
- Schulze, T., Krug, S., and Schader, M. (2012). Workers' Task Choice in Crowdsourcing and Human Computation Markets. *ICIS 2012 Proceedings*.
- Smith, S. (2013). Determining Sample Size: How to Ensure You Get the Correct Sample Size | Qualtrics. Date accessed: 2017-04-19 URL: <https://www.qualtrics.com/blog/determining-sample-size/>.
- Stanford University (2017). Machine Learning | Coursera. Date Accessed: 2017-05-14 URL: <https://www.coursera.org/learn/machine-learning>.
- The Pennsylvania State University (2017). 7.5 - Power and Sample Size Determination for Testing a Population Mean | STAT 500.
- The Scipy community (2017). `scipy.stats.anderson` — SciPy v0.19.0 Reference Guide. Date accessed: 2017-04-21 URL: <https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.anderson.html>.
- Toutin, T. (2004). International Journal of Remote Sensing. *International Journal of Remote Sensing*, 25:10:1893–1924.

-
- von Ahn, L. (2008). Human Computation. In *2008 IEEE 24th International Conference on Data Engineering*, pages 1–2. IEEE.
- Walpole, R. E., Myers, R. H., Myers, S. L., and Ye, K. (2012). *Probability & Statistics*. Pearson Education, ninth edition.
- Wang, X., Goh, D. H.-L., Lim, E.-P., Wei Liang Vu, A., and Chua, A. Y. K. (2017). Examining the Effectiveness of Gamification in Human Computation. *International Journal of Human-Computer Interaction*, pages 1–9.
- Yap, B. W. and Sim, C. H. (2011). Journal of Statistical Computation and Simulation Comparisons of various types of normality tests Comparisons of various types of normality tests. *Journal of Statistical Computation and Simulation*, 81(12):2141–2155.