

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.

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

Contents

- Abstract. v
- Sammendrag vii
- Preface ix
- 1 Methodology and experiment 1
 - 1.1 Web-application 1
 - 1.2 Pilot test 1
 - 1.2.1 Execution of the pilot test 2
 - 1.2.2 Results from the pilot test 2
- 2 Statistics 5
 - 2.1 Sample data from Survey 5
 - 2.2 Statistics theory 5
 - 2.2.1 Normal testing 5
 - 2.2.2 Hypothesis testing. 6
 - 2.3 Survey results 9
- Appendices 11
- A Tets 13

1 | Methodology and experiment

1.1 Web-application

This thesis used an online web-based survey to conduct the experiment. An online survey avoid 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).

In a self selected sample, which is some the case here, there is potentially a bias in the sample (Ben and Plaisant, 2009).

1.2 Pilot test

It is important to pilot test the survey prior to actual use (Ben and Plaisant, 2009). A pilot test provides an opportunity to validate the wording of the tasks, do the participants understand the tasks? It also helps understand the time necessary for completing the survey, which should be communicated to the participants in prior to the survey (Schade, 2015). The pilot-test is conducted with a small sample of users. Results from the pilot-test can be used to determine the sample size. The sample size tells us how many responses that are needed (Smith, 2013). The formula for determining the sample size requires the standard deviation, how much variance to expect in the response (Smith, 2013). This standard deviation can be calculated from the pilot-test results.

A pilot test was conducted with a total of eight participants, five experienced and three non-experienced participants aged from 22 to 64 years. After the pilot test the usability was measured. The standard ISO 9241-11 suggests that measures of usability should cover effectiveness, efficiency and satisfaction (ISO, 1998). Measuring these three classes of metric can vary widely and makes it difficult to make comparisons of usability between different systems. " just because a particular design feature has proved to be very useful in making one system usable does not necessarily mean that it will do so for another system" (Brooke, 1996). Usability in this thesis was measures with the *System Usability Scale*(SUS). This scale gives an subjective measure of usability and was developed by John Brooke. The *System Usability Scale* questionnaire consists of ten statements where the participants rate their agreement in an five-point scale (Ben and Plaisant, 2009). Subjective measure of usability is usually obtained through the use of questionnaire and attitude scales (Brooke, 1996). SUS was developed to be quick and simple, but also reliable enough to be able to compare performance changes between versions (Brooke, 1996).

The usability is important to measure. If the participants doesn't 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. Usability is an important factor to get enough participants to do the whole survey and not quit halfway.

1.2.1 Execution of the pilot test

The pilot test started with a brief information about this study and the survey. They were told to talk out loud during the survey, 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 in the survey correctly. After the survey a *System Usability Scale* questionnaire was answered by the participants. At the end the participants was asked to give general feedback on the web-application. The SUS score and the feedback was then used to determine the usability of the web-application and to determine which improvements to be done.

1.2.2 Results from the pilot test

- Did someone knew micro-tasking? Can we see something here?

The average SUS score was 84.64 out of 100.

All participants thought that the instruction movie was confusing. It was short, the instructions went too fast and it lacked voice descriptions. The movie needed major improvements.

Overall feedback on the tasks was that it was difficult to understand which building was which and also if the building shape layer was selected or not in question one. The lack of labels on the building was done on purpose to get the task as much as possible realistic. The selection of best fitting shape needed improvements, it had to be made clearer that selection was done by clicking on the shape, not by using the layer control as some thought.

Also, some pages had too much information and long sentences. The task progress bar was removed, no one noticed it, only the survey progress bar on the top right is necessary.

The two oldest participants spent almost twice as much time on the test than the younger. Maybe it where 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 where experienced and the other non-experienced, so this is a surprising result. CHECK THE TIME ON EACH TASK FOR THE OLDER PARTICIPANTS.

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.

1.2.2.1 Statistics result

First the pilot-test result need to be normality tested. As mentioned in 2.2.1 this thesis will use the Anderson-Darling test. The python library Scipy has an Anderson-Darling test function (The Scipy community, 2017), this was used to answer the normality test with the Anderson-Darling hypothesis. Testing the total time data form all participants (in total 32 entries) in Anderson-Darling gave a *p-value* of 0.717. The null hypothesis cannot be rejected, then there are two possible cases. One can either accept the null hypothesis or 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 necessary accept H_0 , but fails to reject H_0 (Walpole et al., 2012).

Anderson-Darling test

Data: All participants - total time, 32 entries

P-value: 0.717

P-value > 0.05

H_0 : Accepted, failed to reject

Anderson-Darling test

Data: All participants - total time, 32 entries

P-value: 0.717

P-value > 0.05

2 | Statistics

2.1 Sample data from Survey

Independence of observations. This is mostly a study design issue and, as such, you will need to determine whether you believe it is possible that your observations are not independent based on your study design (e.g., group work/families/etc). A lack of independence of cases has been stated as the most serious assumption to fail. Often, there is little you can do that offers a good solution to this problem.

Designed the survey so that the observations should be random and independent

- Random order on the tests
- Random color on the layers
- Random which order the layers was drawn on the map
- Random which order the metadata was written in the table

2.2 Statistics theory

2.2.1 Normal testing

All parametric statistics assumes normally distributed, independent observations. Parametric tests are preferred in statistics because it got more statistical power than non-parametric 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.

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. 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 2.1, and kurtosis, see figure 2.2, to quantify how far from the normal distribution the sample is from the terms of asymmetry 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).

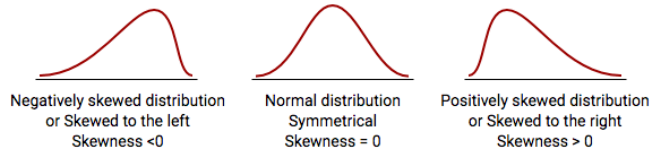


Figure 2.1: Skew (MedCalc Software bvba, 2017)

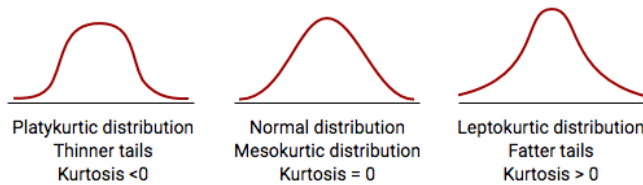


Figure 2.2: Kurtosis (MedCalc Software bvba, 2017)

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 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 necessary accept H_0 , but fails to reject H_0 (Walpole et al., 2012).

2.2.2 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 test which statement is most likely (technically it's 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 first step, after identifying the hypothesis, 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).

2.2.2.1 Two sample t-test

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

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

Then the hypothesis refer to a one-tailed two sampled t-test. 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.

Because the one-sided tests can be backed out from the two-sided tests. (With symmetric distributions one-sided p-value is just half of the two-sided pvalue). It goes on to say that scipy always gives the test statistic as signed. This means that given p and t values from a two-tailed test, you would reject the null hypothesis of a greater-than test when $p/2 < \alpha$ and $t > 0$, and of a less-than test when $p/2 < \alpha$ and $t < 0$.

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

H_0 : Mean task time between participants are equal

H_A : Experienced participants finish the tasks faster, use less time

H_0 : Total number of correct elements between participants are equal

H_A : Experienced participants have a higher number of correct elements

H_0 : There are no difference in total number of correct elements between the tasks

H_A : Participants have more correct elements on the one element task

H_0 : There are no difference in mean time between the tasks

H_A : Participants finish the one element task faster

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

2.2.2.2 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.

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 variable in all groups(Lund Research Ltd, 2013b). The hypothesis test can be written like this:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

H_A : 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, not only the sample (Lund Research Ltd, 2013b). 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.

In an one-way *ANOVA* test there should be one variable and minimum three independent groups, which is an relevant approach considering the data produced from this thesis survey. There are at least two 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* test 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 <i>Variable = time, group = tasks</i>
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 <i>Variable = Number of correct elements, group = tasks</i>

The hypothesis written above will be testen in the section blabla.

2.2.2.3 Wilcoxon Rank-Sum test

The Wilcoxon Rank-Sum test is an appropriate alternative to the two-sample t-test (see subsection 2.2.2.1) when the data are not normally distributed (Walpole et al., 2012). Since this method is nonparametric (or distribution-free) it do not require the assumption of normality. A nonparametric method is much more efficient than the parametric procedure when the set of data used deviates significantly from the normal distribution (Walpole et al., 2012).

2.2.2.4 Kruskal-Wallis test

The Kruskal-Wallis test is a nonparametric alternative to *ANOVA* (see subsection 2.2.2.2) (Walpole et al., 2012). This test should be used if the assumption of normal distribution failed. As mentioned in subsection 2.2.2.3, a nonparametric method does not assume normality.

There are 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).

2.3 Survey results

Appendices

A | Tets

Fbox

Some text esfljsf
lskj lksdjflsk slk

Some text
kduhaszkdh aszkd-
jhs zkjd fh skdj
skd

dwkjdkwjdh wkjdhw kjdh wkjhd qwkjhd kwd qw .

text

dwkjdkwjdh wkjdhw kjdh wkjhd qwkjhd kwd qw .

Bibliography

- Ben, S. and Plaisant, C. (2009). *Designing the User Interface*. Pearson, fifth edition.
- Brooke, J. (1996). *SUS-A quick and dirty usability scale*. "Usability Evaluation In Industry". Taylor & Francis.
- 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>.
- 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.
- ISO (1998). ISO 9241-11:1998(en), Ergonomic requirements for office work with visual display terminals (VDTs) — Part 11: Guidance on usability.
- Lund Research Ltd (2013a). Hypothesis Testing - Significance levels and rejecting or accepting the null hypothesis.
- Lund Research Ltd (2013b). 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>.
- MedCalc Software bvba (2017). Skewness and Kurtosis. Date accessed: 2017-04-23 URL: <https://www.medcalc.org/manual/skewnesskurtosis.php>.
- Motulsky, H. (2013). Intuitive Biostatistics: A Nonmathematical Guide to Statistical Thinking, 3rd edition. In *Intuitive Biostatistics*, chapter 24.
- Pearson, A., Sözcükler, A., düzeltmeli Kolmogorov-Smirnov, L., Pearson ve Jarqua-Bera testleri Derya ÖZTUNA Atilla Halil ELHAN Ersöz TÜCCAR, A., and Öztuna, 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.
- 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/>.
- 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/>.
- The Pennsylvania State University (2017). 7.5 - Power and Sample Size Determination for Testing a Population Mean | STAT 500.

BIBLIOGRAPHY

- 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>.
- Walpole, R. E., Myers, R. H., Myers, S. L., and Ye, K. (2012). *Probability & Statistics*. Pearson Education, ninth edition.
- 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.