Reducing Human Error in Online Controlled Experiments

A case study at ING



Martijn Steenbergen

Reducing Human Error in Online Controlled Experiments

THESIS

submitted in partial fulfillment of the requirements for the degree of

MASTER OF SCIENCE

in

COMPUTER SCIENCE

by

Martijn Steenbergen born in Woerden, the Netherlands



Software Engineering Research Group Department of Software Technology Faculty EEMCS, Delft University of Technology Delft, the Netherlands www.ewi.tudelft.nl



ING Bijlmerdreef 106 Amsterdam, the Netherlands https://www.ing.com/Home.htm

© 2020 Martijn Steenbergen. Cover picture: Photo by Joanna Kosinska on Unsplash.

Reducing Human Error in Online Controlled Experiments

Author:Martijn SteenbergenStudent id:4311442Email:martijn.steenbergen@outlook.com

Abstract

Online controlled experimentation (OCE), also called A/B testing, is an often used tool in industry to determine if deploying changes into production is the right decision to make. Running experiments has shown an immense impact to the revenue of companies in industry, however this type of experimentation comes with a lot of pitfalls, of which some that can invalidate the entire experiment.

This thesis describes the impact these pitfalls have on the work of experimenters at ING, a global bank, by performing informal interviews with practitioners and performing a survey with 52 participants. Next, building on existing solutions, a set of solutions is proposed to solve these pitfalls. To determine if these solutions solve the problem and will help the experimenter, these solutions are validated in the same survey.

This thesis shows that experimenters are well informed about the existence of pitfalls and believe that almost all should be resolved, with the exception of competitor safety, which is believed to not be important. There are many promising solutions to these pitfalls which experimenters rate as helpful. The best rated solution was "Enforcing the correct experiment duration". Almost all respondents perceived the solution to (slightly) help the experimenter in performing their experiments. Finally, this thesis creates a roadmap for evaluating these solutions in a real-world scenario.

Thesis Committee:

Chair:Prof. Dr. A. van Deursen, Faculty EEMCS, TU DelftUniversity supervisor:Ass. Prof. Sebastian Proksch, TU DelftCompany supervisor:Kevin Anderson, INGCommittee Member:Ass. Prof. Jan S. Rellermeyer, Faculty EEMCS, TU Delft

Preface

Dear wonderful reader,

Thank you for being awesome and reading this thesis. 5 years ago, I started my journey at the TU Delft. Never could I have imagined that I was starting the most amazing adventure of my life. Your kindness and knowledge has made me the person I am today. Thank you for that.

They say that a master thesis is the first thing you really do on your own. Although that is partially true, there are so many people who have supported me and without them, this thesis wouldn't be what you find here today.

Thank you to my supervisors, Arie, for giving me the freedom to do what I like while always pointing me in the right direction. Thank you, Sebastian, for pushing me to see the facts and helping me create something actually achievable. Thank you, Kevin, for bringing me into contact with all of the amazing people at ING and for being there. It was an amazing opportunity to see you work and shine.

Thank you, all the people that proofread my thesis. Thank you, Irene, for always reading my thesis and giving me feedback (even when it was bad) and helping me plan in the beginning days. You are an amazing addition to the team. Thank you Chiel, for the mountain of feedback that I hope I did proud.

Thank you, Team Tetris for welcoming me with open arms, letting me work in your codebase and helping me with all the problems I had. Thank you, survey respondents for making my research possible. Thank you, Huiskool, for helping me grow into the person I am today. You are always there for me. Thank you, 4th floor, for making these last years amazing. Thank you, everyone from the Teaching Team, SERG, Flying Fish and so many others that I met during these years. The greatest honor I could have is that I gave you as much happiness as you gave to me.

Thank you.

Contents

Pr	eface		iii
Co	ontent	S	v
Li	st of F	ïgures	ix
1	Intro	oduction	1
	1.1	Online Controlled Experiments	1
	1.2	Human Error in Controlled Experiments	2
	1.3	Industrial Context: ING	3
	1.4	Research Questions	3
	1.5	Contributions	4
	1.6	Thesis outline	4
2	The	context of ING	5
	2.1	Experimentation at ING	5
	2.2	The Experiment Office at ING	6
	2.3	Interviews	11
3	Com	mon pitfalls in Online Controlled Experiments	15
	3.1	Falsifiable hypothesis	15
	3.2	Direction of change in hypothesis	17
	3.3	Guardrail metrics	17
	3.4	Technical debt	18
	3.5	Competitor safety	18
	3.6	Churning users	18
	3.7	Minimum effect size	18
	3.8	Minimum duration	19
	3.9	Withholding results	19
	3.10	Simultaneous experiments	19
	3.11	Failure checks	20

CONTENTS

	3.12 Number of changes	20
	3.13 Day of week effect	20
	3.14 Not confirming the winning variant	21
	3.15 Encourage more experiments	21
	3.16 Higher level question	21
	3.17 Share learnings	22
	3.18 Rerun experiment when results are marginal	22
	3.19 Validation of experiment	22
4	Solutions to pitfalls	23
	4.1 Existing AB-test frameworks	23
	4.2 Solutions	24
5	Research methods used	35
-	5.1 Survey contents	35
	5.2 Participation	36
	1	
6	Results	39
	6.1 Overview of results	39
	6.2 Differences among pitfalls and their solutions	41
	6.3 Pitfall results highlights	44
	6.4 Results per solution	45
7	Methodology for evaluating solutions in a real world scenario	49
	7.1 Research methods used	49
	7.2 WebTrekk	50
	7.3 The platform survey	52
8	Conclusion & Discussion	57
		51
	8.1 Conclusion	57
	8.1 Conclusion	57 58
	 8.1 Conclusion	57 58 59
	 8.1 Conclusion	57 58 59 61
	 8.1 Conclusion	57 58 59 61 61
Bi	 8.1 Conclusion	57 58 59 61 61 63
Bi	 8.1 Conclusion	57 58 59 61 61 63 69
Bi A B	 8.1 Conclusion	57 58 59 61 61 63 69 74
Bi A B	 8.1 Conclusion	57 58 59 61 61 63 69 74
Bi A B C	 8.1 Conclusion	57 58 59 61 61 61 63 69 74 77
Bi A B C	 8.1 Conclusion	57 58 59 61 61 63 69 74 77 77
Bi A B C D	 8.1 Conclusion	 57 58 59 61 61 63 69 74 77 77 87

D.2	Mann-Whitney U test																							9	0
-----	---------------------	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	---	---

List of Figures

2.1	Step 1 of the experiment office flow, picking the experiment type and choosing where the experiment will be run	8
2.2	Step 2 of the experiment office flow, creating documentation for the experiment (title, hypothesis and variants)	9
2.3	Step 3 of the experiment office flow, choosing the participants group and run- time of experiment	10
2.4	Overview of the number of experiments run at ING	10
2.5	Overview of the number of experiments run in production at ING	11
2.6	Number of experiments ING squads performed in the past two years. Each color is a squad. Names have been left out for non disclosure reasons	12
4.1	Screenshot of the second step of the experiment creation workflow	26
4.2	Screenshot of the Google Firebase website. Screenshot taken from https://co	
	<pre>nsole.firebase.google.com/u/0/project/fir-demo-project/config/</pre>	
	experiment/results/3?hl=en	28
4.3 4.4	Screenshot of the metrics dashboard from Airbnb	29
	formed by Meeuwsen et al.	29
4.5	Sample size calculator from Optimizely	30
4.6	Sample size calculator from Adobe Target	31
4.7	Screenshot of the Google Firebase website. Screenshot taken from https://console.firebase.google.com/u/0/project/fir-demo-project/config/	
	experiment/results/3?hl=en	32
4.8	Preview of the experiment given by Optimizely	34
6.1	Participants years of experience	40
6.2	Percentage of participants who experienced the pitfall	41
6.3	Percentage of participants who is aware of the pitfall	42
6.4	Divergent stacked bar chart [12] of the perceived severity of the pitfall	42
6.5	Percentage of participants who believes the solution solves the pitfall	43
6.6	Divergent stacked bar chart $[12]$ of the perceived usefullness of the solution \ldots	43

LIST OF FIGURES

 7.1 7.2 7.3 7.4 7.5 7.6 	Implementation of the full week effect solution	50 51 53 54 55
7.6	Reminder email Diagram showing the combined number of responses for pitfalls, set to the years of experience	56 60
A.1 A.2	Mockup of a possible version 2 of the experiment office	70 71
C.1 C.2	The first page of the survey	78 79

Chapter 1

Introduction

In 2014, Facebook engineers came back after Christmas and found that they had an problem. During Christmas, people had uploaded an enormous amount of Christmas pictures (more than Flickr had gotten in the entire lifetime). This also meant that the amount of content reported as hate speech, sexual abuse, etc. went up so drastically the team responsible for handling reports couldn't keep up anymore. When Facebook engineers started to look at the reports, they realized that 97% of these reports were miscategorized. Weirdly, in many cases the people making the report were in the picture. In fact, a picture flagged as hate speech was usually nothing more than an embarrassing image. To solve this problem, the team made it their goal to reduce the amount of reports by having the person contact the poster of the image to take it down. To achieve this, they added "Its embarassing" to the possible reasons of reporting an image. When a user clicked that it would open a new conversation with that person. The team wanted to run an experiment to test what would work better: an empty chatbox (version A, the control) or a chatbox with a default message (version B, the treatment) [30]. The result was that with an empty chatbox only 20% of the people would start a conversation, but with the default message it was 50%. [43]. This default message was implemented for a long time on the site. This type of test is called an Online Controlled Experiment [30]. This type of experimentation is now used throughout the web and is used extensively by Facebook. For any user of Facebook, there is an almost 100% certainty they have been part of a controlled experiment. In 2015, a user was on average part of 10 experiments at any given time. [43]

1.1 Online Controlled Experiments

Online controlled experiments (OCEs) or in literature also called A/B test, randomized experiment (single-factor or factorial designs), split test, Control/Treatement test and parallel flight are everywhere on the internet [29, 30]. An online controlled experiment is where all visitor to a webpage are split between two different variants to see which performs better on some metric. The idea of a controlled experiments is relatively simple and dates back to Sir Ronald A. Fisher', a founding father of statistics [23]. The idea of systematically running experiments to find the return-on-investment of new features in software seems to

have been first termed by Dan McKinley in 2012 in a presentation about experimentation at Etsy [10]. This idea took off, as the web provides an unprecedented opportunity to evaluate ideas quickly using controlled experiments [30]; deploying changes on the internet is cheap and each experiment can exposes hundreds of thousands —and sometimes even tens of millions— of users to a new feature or change [24, 30]. Furthermore, controlled experiments enables the experimenter to control for external factors, making sure the change is related to a new feature, not to a random fluctuation [36]. Because of this controlled experiments are becoming a standard in data-driven software companies. [19, 18]

Controlled experiments can bring a form of clarity to an organisation. Every experiment needs to have an Overall Evaluation Criterion (OEC), which is a quantitative measure of the experiments objective [13]. These objectives should relate to the overall goal of the organization. This aligns the entire company to focus on one or more (ranked) goals [19, 30]. In the example in the introduction, this goal was number of started conversations. If number of sales is the goal, increasing prices to make more profit, will likely be bad idea, as this will lower the amount of sales.

1.2 Human Error in Controlled Experiments

The number of companies using experimentation and the number of experiments inside these companies are growing [36, 15]. Therefore it is infeasible to have only experts doing experiments. Furthermore, finding an expert in both experimentation *and* the context of the experiment is even less feasible.

To solve this, organizations use in-house tools [18] to help non-experts create and execute experiments. However, this tooling needs to embody the knowledge of an expert. This has proven difficult, as experts are still uncovering and solving problems related to OCEs. In every area, at every stage of running experiments pitfalls exist that can influence the process, at worst invalidating the experiments. Bad data can be actively worse than no data [19]. These pitfalls range from making sure the tooling to collect data are working correctly and making sure the platform is statistically valid, to making sure that the experiment layout is valid. A large amount of literature details these pitfalls and how they have affected the organization reporting it [19, 30, 10, 22].

We can define two main categories of pitfalls, those on a platform level and those on an experimenter level. On a platform level it is a must to make sure that the data is gathered correctly. This problem is, although important, only interesting in a technical perspective. However, pitfalls on an experimenter level combines the technical with a human aspect. That is why it was chosen focus on pitfalls, which contain *human error*.

Many companies are moving to become data-driven organizations and therefore are starting to use experimentation to make decisions. When building their own platform, a list of implementation guidelines is missing. This thesis will provide indications which areas might provide the biggest payoff.

1.3 Industrial Context: ING

ING is a global bank with its main base in the Netherlands and Belgium. It employs around 53,000 employees, which serve around 38.4 million customers, corporate clients and financial institutions in over 40 countries [27]. ING is already running some controlled experiments on the web and wants to expand this in other areas of the company. It can therefore be determined that it is in the midst of adopting this technology. ING has their own experiment platform called the Experiment Office, where tests can be created and results be viewed. ING has recently started a partnership with the TU Delft under the banner of AI4FinTech, under which this research was performed.

1.4 Research Questions

From the main problem statement follows **the main research question** this thesis tries to solve:

How can we reduce human error in online controlled experimentation?

To help experimenters, many areas will need to be resolved. Therefore this main question is divided in 3 sub questions. First, an overview needs to be created of the pitfalls, which are relevant for this research and solutions which can solve these problems. Therefore the first sub research question is:

RQ1 What are key pitfalls and their potential solutions?

After a list of solutions exist, they need to be validated. As one important aspect of supporting the experimenter is that any change we make should be in the benefit of the experimenter. Any change should not decrease the satisfaction the experimenter has of the platform. This requirements is both important, because it is the right thing to do and because organizations have an incentive to keep this in mind. After all, when changing to become a data-driven organization it is counter productive to have tools employees/experimenters do not like. Furthermore, a deeper dive into the pitfalls is required, as the same holds for the pitfalls themselves. Having a great solution does not mean something if the pitfall in and of itself is not believed to be important.

RQ2 How are pitfalls and their solutions experienced?

Once a clear idea exists of which solutions could work, implementing such solutions in a real-world environment should be the next priority. A plan needs to be layed out on how to implement the solution, solve problems on the way to implementation and and if the results after implementation collaborate the results of earlier work. Therefore the last subquestion is:

RQ3 How can one evaluate a solution in a real-world scenario?

1.5 Contributions

This work made many contributions to the scientific knowledge of continuous experimentation. To the authors knowledge, this is the first body of knowledge surveying experimenters with solutions for pitfalls, before implementation. In the rest of this section these contributions will be discussed.

First, **this thesis proposes solutions** to pitfalls in online controlled experimentation by looking at existing solutions.

Second, it **shows a comprehensive overview of the state of online controlled experimentation at ING**. By performing an interview with practisioners in the field, determining workflows and documenting experiences by practisioners and analyzing the current experiment platform, an empirical contribution[6] is added to the knowledge of OCE's.

Finally, another an empirical contribution[6] is created, as this thesis evaluates how 19 pitfalls are experienced and verifies the 8 proposed solutions, which can help experimenters avoid those pitfalls when experimenting.

1.6 Thesis outline

To have a good understanding of the industrial context (ING), in chapter 2 a small case study is done, to give the reader an overview of the state and practice of online experimentation at ING. This chapter concludes with informal interviews conducted with ING employees to get an overview of their workflow (2.3).

In chapter 3, the first part of RQ1 is answered; which pitfalls in existing research are key to this thesis. Each of these key pitfalls is then explained to the reader in greater detail. In chapter 4, the second part of RQ1 is answered: which potential solutions exist to solve the pitfalls. In chapter 5, the methodology is described on how a survey is used to determine how pitfalls and their solutions are experienced by experimenters (RQ2). After this survey has been performed, in chapter 6, the results of the survey are analyzed. These results can then inform the implementation of the solutions. In chapter 7, the thesis proposes a roadmap to evaluate the proposed solutions, in order to answer RQ3.

Finally, in chapter 8 the thesis is wrapped up, by drawing conclusions, discussing threats to validity and describing future work.

Chapter 2

The context of ING

Before we can dive deeper into the problem, we need to first have an overview of the context in which this problem resides.

The goal of this chapter is to present the reader a thorough view of the state of continuous experimentation at ING. This is done by first in section 2.1, by providing the history of experimention at ING and the place in the run walk and fly stage. Next, in section 2.2 the experiment platform of ING is discussed. Finally, in section 2.3 the interviews that were performed to gain a deeper understanding of the state of continuous experimentation at the company.

2.1 Experimentation at ING

ING is a global bank with its main base in the Netherlands and Belgium. It employs around 53,000 employees, which serve around 38.4 million customers, corporate clients and financial institutions in over 40 countries [27]. ING, just like many other companies is in the midst of digital transformation [25, 26]. More than 80 percent of customer interactions now go through mobile devices. This switch from offices to digital devices is so profound that the main message in the yearly report of the CEO at the time, Ralph Hamers, was that "The digital customer experience is the key differentiator". He explains that "ING's ambition is to be a leader in terms of the digital banking experience, offering retail and wholesale customers everywhere the same empowering and differentiating experience" [26]. Therefore optimizing this customer experience should be at the heart of the future of ING. One of the ways in which ING has been doing this, is via ING's Analytics Unit. Established in 2018 [26], its goal is to accelerate the bank's analytics capabilities and lead the transformation to become a truly data-driven company [26]. Currently more than a 100 data scientist are working at this unit to make this goal a reality [26].

History of Experimentation at ING From our conversations with employees, we learned that ING has a long history with experimentation. At the PostBank, a predecessor of ING, direct mail was sent to a subset of customers. The direct mail consisted of multiple parts, a header, images, multiple paragraphs, etc. For each of these parts, variants were created and

combinations of these were sent to customers. The most successful variants were combined into the final mail, which was sent to all customers.

2.1.1 Maturity of continuous experimentation at ING

To determine the state of ING testing, we will be using the *crawl walk run and fly* stage by Fabijan et al. [17]. The authors are experts in the field of continuous experimentation, however the framework that is described here does not seem to be referenced in in literature by other authors. Therefore we can consider this framework as highly indicative, but not conclusive to the state of experimentation at ING. According to this framework, ING is in the walk stage, similarly to most other companies trying to perfrom OCEs [18].

ING has made great progress in implementing Continuous experimentation, however there are still improvements that can be made. ING and the research community could use this to their advantage. ING is in the perfect position to document the struggles, roadblocks and solutions, best practices and strategies to lead other organizations moving to a data driven culture.

2.2 The Experiment Office at ING

Experimentation at ING is done using the Experiment Office. The Experiment Office is a tool built in-house to support their own experiments.

The Experiment Office is decoupled from the implementation of the experiment. This is done to support multiple platforms at once. Currently the web is the only supported platform on which experiments can be run. However ING is slowly rolling out support to enable running experiments on the mobile apps as well.

The Experiment Office accomplishes its decoupling, by having a web API, which a developer can call to register feedback. By calling endpoints on the web API the server can determine which variant this user came from and if the variant was successful. The Experiment Offices also allows the creation of experiments in a test or acceptance environment to test if the test is correctly set up.

One interesting thing to note is that the Experiment Office is unable to determine if a variant failed for a user. It's only able to determine that a user succeeded. There is no way for the implementation to register a failed instance and should be inferred from other metrics. This is a result of how the platform was implemented.

When creating an experiment, an experimenter might want to test the experiment in a test or acceptance environment before deploying to production. The Experiment Office offers this functionality. By having each environment on a different subdomain, a clear division is made to experimenters on where they are running their experiments.

The experiment is maintained by the Tetris team at ING. This team took over the Experiment Office from squad Panama in March of 2018. With a small team, they maintain and extend the Experiment Office.

2.2.1 Types of experiments

The Experiment Office has multiple ways of showing users different variants. Although only A/B-tests can be seen as an actual experiment, we will follow the documentation of the Experiment Office:

- **A/B test** An A/B-test is an online controlled experiment, as follows from the definition. Within ING, the term A/B test is preferred over "online controlled experiment".
- **Pilot** A pilot differs from an A/B test in the participation strategy. Instead of using a random subset of users to show variant A or B to, the experimenter can either provide a list of users to show the new variant to or users themselves can choose to see the new variant. This therefore introduces bias into the experiment and cannot be used as an actual experiment. However a Pilot can still be very useful to gather feedback from active users. Or to make sure a feature actually works before rolling it out.
- **Rollout** To reduce the possibility of shipping a broken feature to customers, a rollout experiment type exists that enables the experimenter to slowly roll out a new feature to all users. As soon as it is manually determined that the future is not working as expected, if the new feature is breaking functionality or it is throwing exceptions, the feature can be rolled back without having been seen by the entire user base of the platform.

2.2.2 Creating an experiment

To create an experiment, two main steps are required. The first step is to setup the new experiment in the Experiment Office. The second step is to actually build experiments where the content exists (Polymer [21], CMS), commit this to the experimenters own repository and deploy this. To create new experiments in the web interface of the Experiment Office, an experimenter navigates to the dashboard and clicks new experiment. There they can choose the experiment type, as seen in Figure 2.1, which was discussed above.

After this is done, they can move to step 2 (Figure 2.2), where more information about the experiment is requested: the title and a hypothesis. This forces experimenters to think about the experiment and what they actually want to test and what the result of the experiment should be. Next up they choose a success criteria where they are forced down to decide when this experiment should be considered a success.

Finally, in the last step of the process, the experimenter chooses the target group and the tool/tracking method used to collect metrics and track the user through the system.

2.2.3 Experiments on the platform

The Experiment Office launched in March of 2015. The first experiment in production launched a year later, in December of 2016. A total of 1875 experiments have been run on the platform of which 1067 were an online controlled experiment. To get a better view of the experiments being run at ING, an export was made of the platform at June 26th 2020. The rest of this section will discuss insights from this data.

2. The context of ING

		🔯 Martijn (M.J.W.) Steenbergen	Give your feedback Logout
ING 🍌 The Guide	Fruitloops A/B Testing		
Dashboard Experiment Over	view Experiments Squad Ranking Sea	Documentation	
> Home > Experiments > Create expe	riment		
Experiment Conf	igurator		
1 Setup	2 The test	3 Target group	\sim
E			
Experiment type	A/B test		
	○ Rollout		
Implementation Pattern	○ Content A/B test on ING.nl		
	 Feature A/B test (for developers) 		
	 Content A/B test on Orange Channels (I 	D20)	
O You are now in the Produ	tion environment.		
The experiments that you If you have not done it ye	create here are not available in Test or Acceptan t we recommend to start building and testing up	ice. ur experiments from there	
The respective dashboard	s can be found here: <u>Test</u> , <u>Acceptance</u>		
	Next >		

Figure 2.1: Step 1 of the experiment office flow, picking the experiment type and choosing where the experiment will be run

If one looks at the A/B tests performed from June 26th 2018 until June 26th 2020, 339 experiments were started, 43% in production, 37% in testing and 20% in acceptance.

From Figure 2.4 we can see that until March 2019 the number of experiments increased, however after this time, the number of experiments are slowly going down. If we look only at experiments in production (Figure 2.5), this change is even more pronounced. We can also see a sharp decline in the number of A/B experiments in the months of March, April, May and June 2020, which can probably be attributed to the corona outbreak of 2020 [49].

As there are three environments, a single experiment-id might exist on more than one environment. In the Experiment Office, 92 experiments are found in two stages, 148 in all three.

Many experiments do not yield significant results. If we look at A/B experiments run in production, only 26% (37 experiments) of the experiments are significant. In test and acceptance, only a single experiment is significant, which is explained by the low traffic on these pages, making it hard to reach significance. Many teams or squads, as they are called within ING, are using the platform, as can be seen in Figure 2.6, which is both a positive and

ashboard Experime	nt Overview	Experi	iments Squad Ranking Search Documentation	
kperiment C 1 ^{Setup}	onfigur	ato	2 The test 3 Target group	>
Title and key	Experiment Experiment	title key	0	
	Feature		Value is used by the developers. Valid characters are: "a-z0-9"	
Hypothesis 🕕	Template By will lead to measured by	Free-1	text replacing a drop-down with a category page sales increasing more orders	
/ariants Control	Key Description	(i) (i)	control-group	
Variant	Key Description	() ()	variant1	
			+ Add variant	

Figure 2.2: Step 2 of the experiment office flow, creating documentation for the experiment (title, hypothesis and variants)

2. The context of ING

Home >Experiments >Create expe	riment		
xperiment Conf	igurator		
1 Setup	2 The test	3 Target group	
Audience type	 Customers (ACMA id) 		
	 Customers (RGB) 		
	Employees (corporate key)		
Choose target audience	All customers		
	Customers who enter with a certain referer heade	r	
	 Specific customers 		
Size of group	Unlimited •		
Start date	28-12-2018		
End date	11-01-2019		

Figure 2.3: Step 3 of the experiment office flow, choosing the participants group and runtime of experiment



Figure 2.4: Overview of the number of experiments run at ING



Figure 2.5: Overview of the number of experiments run in production at ING

negative result. It is positive, as this means it is likely many parts of the organisations are learning to become data-driven. However, at the maturity level (2.1.1) ING is in now, this could lead to knowledge being spread out too thin and therefore experimenters not having all the knowledge to run trustworthy experiments and get trustworthy results.

2.3 Interviews

Overview To get an overview of the difficulties experimenters face when creating controlled experiments at ING, we conducted interviews with 4 users (experimenters) (P1-P4) of the Experiment Office. These interviews were unstructured and performed in an informal way with the goal to get an overview of how online controlled experimentation is performed at ING. Questions included how they come up with an experiment, how they use the Experiment Office to run their experiment, what they think of the Experiment Office and what can be improved. Because the interviews were performed by a single person, and in an informal way, this can introduce bias into the results, both in the way the questions were asked and how the results are presented here.

Participants Participants were chosen to have many different viewpoints (experience, platform, etc.). Participants were requested via mail to have an interview of around 45 minutes. Each of the participants at the time of the interview had run experiments in the last month on the Experiment Office. The number of conducted experiments varies per participant. Numbers ranges from 5 (P1), 10 (P4) or 14 (P3). The experience of the participants differs from six months (P1) to 9 years (P2). P1, P2 and P3 all use a content management system (CMS) where they create two different versions of the same page and hook the Experiment Office into it. P4 is a programmer, who uses code to create experiments and



Figure 2.6: Number of experiments ING squads performed in the past two years. Each color is a squad. Names have been left out for non disclosure reasons

finds that creating experiments in this way takes a long time. All see the potential of using experiments to improve customer experience.

Other areas in OCE Interesting to point out is that experiments at ING are not only performed using the Experiment Office, but also using the tools from social media platforms, where ING puts advertisements. For example, Google ads or Facebook ads are used to serve potential customers ads for new products. P2 mentioned that there might be more experiments performed on these platforms than experiments on the ING.nl website. In this case the titles of the ads are changed or it is tested if a countdown timer or static time works best (P2).

Another area where ING performs experiments is via email. If a customer is already an ING customer, ING is allowed to contact the customer to sell a new product. According to P2, the moment this email is sent is crucial. P3 mentions the journey a customer takes when buying a product; See (realize you want something to solve a problem), Think (researching the best solution), Do (buy the product), Care (evaluating if they bought the best solution). Only in the Think and Do phases an email can be relevant and this poses the problem of figuring out in which stage the customer is and how one can best help them.

Origins of OCEs For different participants, inspiration for new experiments comes from different angles. P3 goes through the process the customer goes through and notes down everything that makes P3 think: "hmmm, that can be improved". Then customer feedback is added to add new ideas and determine which ideas have the largest chance of being impactful. P3 states that this order is a must, as you need context to understand what the customer is going through, before reading the feedback. P2 has a different approach and estimates that for them the feedback only accounts for around 20% of new experiments. For P2 the largest origin (~65%) of the new ideas for experiments come from conversations with other squads. Lessons are learned in one squad, which can be relevant for other parts of the website. The final 15% are found in customer feedback, where they complain that something does not work correctly or have other ideas on improving the flow.

Creating an OCE After ideas have been brainstormed, the experiment needs to be launched. P3 mentions that they calculate the minimum duration of an experiment up front, although both P2 and P3 mention that any experiment will be run for at least two weeks. An observation, made explicit by P2, but also applicable P1 and P3, is that once an experiment is running, another one cannot be launched. This is impossible on the platform with CMS experiments and generally considered bad practice. This creates the situation where if there are 26 small improvements you would like to make to a page, this would take an entire year. This is why P3 mentions that 3 weeks is the maximum length an experiment can run, as otherwise it would simply take too much time, even if the results are not significant yet. This is also recognized in the literature [19].

Dissapointing outcomes In the end, conclusions need to be drawn. Something that is mentioned in literature [17] and here is reiterated, is that a negative outcome can be disappointing. P3 mentions that at ING, the best practice of leaving the current version in production is followed, after which something else might be tried. However the experimenter will always stay defensiveness about the idea. They might think of possible ways that their idea might be right after all. This enforces the idea that trust in the platform is of upmost importance.

Drawing conclusions It is also interesting to see that for drawing the conclusions multiple approaches are taken. P1 and P4 look at the entire sales funnel to draw the final conclusion. A sales funnel is the entire process from seeing the ad to buying the product. After all, more people clicking on an advertisement does not necessarily imply more people buying the product. If more people click on the advertisement, to continue the example, P1 will put it into production and check after two weeks if the number of sales has also increased. P2, however, states that they do not check the sales funnel, as there are too many factors involved to draw good conclusions. Before making a decision, P4 also looks at some guardrail metrics (3.3) to make sure nothing is breaking.

Thoughts of the experiment office Overall, participants are pretty happy with the Experiment Office and praise its ease of use. P1 found getting started with the tool easy. The par-

ticipants were also asked if there could be anything improved about the Experiment Office. Interestingly, all participants mentioned that they would like to keep track of more metrics. P2 and P3 want more metrics to base their conclusion on. Besides this, P4 would also like to have guardrail metrics. P1 would even like to follow the customer through the entire sales funnel and mentions that the integration of the metrics software, WebTrekk (7.2), could be improved upon. P1 and P4 mention that grouping is also lacking, for example P4 would like to group by device (who is on mobile).

Another point of improvement is the sharing of results. Both P1 and P2 miss the possibility to share knowledge about which changes have the biggest effect on the KPI's (key performance indicators). Especially since P2 mentions that such a large part for new ideas for experiments comes from other squads.

When talking about how P3 calculates the minimum duration, they mention they are using an Excel file, where they fill in fields and it gives them the length of an experiment. They would rather see this as a part of the Experiment Office.

P4 mentions that the Experiment Office does not maintain any libraries to integrate with code (in later meetings with the team that maintains the Experiment Office, we learned that they do not have the manpower to do so). P4 continues explaining that because of this, many teams, including his, have created their own libraries to interact with the platform, making the cost and buy-in for experimentation higher and harder than necessary.

Conclusion This sections paints a picture of how experiments are being performed, from the origin of the idea for testing, to drawing the conclusions. Interesting points are that ING is performing experiments, not only on web, but also on other platforms, like Facebook and Google ads. There are still many features that experimenters want, but overall the opinion of the Experiment Office is positive.

Chapter 3

Common pitfalls in Online Controlled Experiments

This chapter will be based on work of Mulders [40], who created a reprehensive literature review of 35 pitfalls found in the field of online controlled experimentation. These pitfalls are often made mistakes, while performing online controlled experiments. Other attempts to find additional pitfalls yielded no new results.

However, not all of these pitfalls are relevant for this research. Some of the pitfalls are already resolved, meaning ING has already implemented solutions to these problems. Trying to figure out the impact or possible solutions is no longer relevant in this case. For example, any experiment started at ING must already have a hypothesis. Asking how often they create an experiment without hypothesis, dives more into the effectiveness of the solution and although interesting, is out of scope of this thesis.

Other pitfalls do not contain a human element. As described by the research questions (1), any pitfall not caused by human error merely requires implementation effort and not are therefore requires a different (less interesting) way to be researched. For example, checking for data quality issues is not something that requires the interaction of an experimenter.

To make sure the filtering of the pitfalls was done correctly, the results were verified by an ING employee, who has extensive experience with OCE's.

Table 3.1 shows that 19 pitfalls are both human error and not resolved in the platform of ING. These we will explain further down below.

Something to mention is that there are three pitfalls, which Ernst [40] does not provide a reference in the original work nor for which references could be found. These pitfalls are "Withholding results", "Not confirming the winning variant" and "Encourage more experiments".

3.1 Falsifiable hypothesis

Issue There is no hypothesis for the experiment or the hypothesis cannot be proven wrong [19]

Id	Shorthand	Is Human Error	Is Implemented
01a	Hypothesis exists	Yes	Yes
01b	Falsifiable hypothesis	Yes	No
02a	Available metrics	No	No
02b	Direction of change in hypothesis	Yes	No
02c	Guardrail metrics	Yes	No
03a	Technical debt	Yes	No
03b	Competitor safety	Yes	No
03c	Churning users	Yes	No
04a	Minimum effect size	Yes	No
04b	Minimum duration	Yes	No
04c	Withholding results	Yes	No
05a	Simultaneous experiments	Yes	No
06a	Failure checks	Yes	No
07a	Needed metrics are present	No	No
07b	Dedicated metric collection system	No	Yes
08a	A/A tests	No	No
08b	P-values in A/A tests	No	No
08c	Significant A/A tests	No	No
09a	Number of changes	Yes	No
10a	Cross experiment contamination	No	Yes
11a	Intermediate conclusive results	Yes	Yes
11b	Day of week effect	Yes	No
11c	Early stopping	No	No
12a	Post experiment health checks	No	No
12b	Confirm winning variant	Yes	No
12c	Post experiment checking of guardrail metrics	No	No
13a	Encourage more experiments	Yes	No
14a	Grouping of experiments	No	No
14b	Higher level question	Yes	No
15a	Share learnings	Yes	No
17a	Data quality issues	No	No
18a	Novelty effects	No	No
19a	Skewed data	No	No
20a	Rerun experiment when results are marginal	Yes	No
20b	Validation of experiment	Yes	No

Table 3.1: table of all the pitfalls . The ids are the ids as they originally appear in the work of Mulders [40]

- **Cause** When creating an experiment, it is of great importance to know what the goal of the experiment is going to be. One important part of determining that goal is the hypothesis that describes what the experiment is trying to prove.
- **Result** If a hypothesis cannot be proven wrong, there will be no experiment necessary, as the result is already known.

3.2 Direction of change in hypothesis

- **Issue** It is unclear if the experimenter is looking for any change (better or worse) or only wants to know if the metric moves in one direction (only better). An example of such a correct hypothesis is: "When removing the sign up, we expect signups to go down". In this case the metric measured is the number of signups and the direction is that this number will decrease.
- **Cause** Statistically, the maths between calculating significance and calculating the possibility of a significant change in one direction is different than calculating the possibility of a significant change in one of two directions. Therefore, by forcing the experimenter to choose, we can better manage this. Another advantage of forcing the experimenter to do this, is that they are forced to think about what the result of the experiment should be [28].
- **Result** If the experimenter does not specify what type of statistical test they want to do, this could result in an invalid experiment and therefore wrong results.

3.3 Guardrail metrics

- **Issue** When running experiments, it is best practice to have a set of metrics setup which are made sure to not change significantly negatively [34, 14]
- **Cause** OCEs can be seen as a playground. A safe place to experiment with new ideas, without facing repercussions when things go wrong. However, just like a parent in a playground, you want to keep an eye out, just in case something dangerous happens. A similar thing occurs with OCEs. When running experiments and following the best practice of focussing on only one metric [30, 14], it is best practice to have some metrics being tracked in the background to make sure nothing is going wrong. The most simple thing to look out for is if an experiment crashes the entire website (see section 3.11). However, there might be other metrics the company wants to measure. For example, it might not matter that more people click a link, if that amounts to people overall creating less revenue. Other metrics which could be tracked by the company are long term goals it is striving for [19]. The collection of these metrics is called guardrail metrics.

Problem If no guardrail metrics are setup, some experimenter might inadvertently ruin the reputation of the entire company or damage the company in another way while focussing on that metric

3.4 Technical debt

- **Issue** Technical debt is the implied cost of additional rework in the future by choosing an easy (limited) solution now instead of taking a better approach that would take longer [47].
- **Cause** One thing that is likely is to happen is leftover code of past experiments being in the code base [44].
- **Problem** This should be removed from the code base to reduce technical debt [19], as these 'shortcuts' can hinder future progress. Therefore making it take longer to implement new features.

3.5 Competitor safety

- **Issue** Running OCEs in production might give a competitor an indication of a new product (line) being developed
- **Cause** Testing product ideas in production causes customers and possibly competitors to become aware of these ideas
- **Problem** This might give competitors more time to catch up or beat time to market [19]. This situation is unfavorable and therefore should be thought about when starting an experiment.

3.6 Churning users

- **Issue** Churning in this context explains the situation that the software is degraded in such a way that users leave the product.
- **Cause** Although it is great to test every idea, experimenters need to always keep the final experience for the users of the platform in mind.
- **Problem** If, for example, an OCE gives customers a bad onboarding experience, they might not choose to join after all, making the organization lose a customer [19].

3.7 Minimum effect size

Issue The minimum effect size is part of the calculation of how long the experiment should run.

- **Cause** The minimum effect size enables one to determine, with statistics, how many users are needed to test with to measure a 5%, 10% or 100% change [19]. After all, it is easier to see if a new version is 100% better than 1%, as the cause of a 1% change is more likely randomness if there is a small number of participants.
- **Problem** If this is not calculated, the experiment could run longer than necessary or too short to get meaningful results.

3.8 Minimum duration

- **Issue** The minimum duration is how long an experiment should at least run to gain significance.
- **Cause** A lot of other pitfalls fall under this main pitfall [14]. Novelty effects, which is the effect of the interaction with a new feature being higher, because the feature is new and exciting [19]. After some time, this newness fades and actual results can be drawn. The reverse can also happen. Changing the layout of a webpage, for example, might decrease the speed and interaction at the beginning, but might perform better in the long run. This effect is called the primacy effect. Another pitfall related to the minimum duration is stopping the experiment when significance has been achieved [36]. Running the experiment longer might result in the feature not being significant after all, as the significance can change over time. [19]
- **Problem** If this is not calculated, the experiment could run longer than necessary or too short to get actual results.

3.9 Withholding results

- **Issue** Withholding results from their experiments from experimentersmight will result in the experimenter drawing better conclusions.
- **Cause** If intermediate results are visible for experimenters and variant A starts out as better than variant B, this might lead to "rooting" for this version to win.
- **Problem** This makes the experimenter biased towards a solution. This should be avoided as much as possible.

3.10 Simultaneous experiments

Issue Experiments which influence each other.

Cause Experiments running at the same time might influence each other. The most extreme example is when one test tests a different color for a button, while another test removes the button.

Problem If the participants are not well distributed, this could lead to unwanted consequences. Furthermore, if users are being exposed to multiple AB-tests at once, this could lead to invalid results [31].

	Color A	Color B	Total
Shown	10	30	40
Hidden	30	10	40
Total	40	40	

Table 3.2: Table showing an example of simulatenous experiments

In Table 3.2, both experiments and variants have an equal number of participants. However, for 75% of the participants the variant with Color A is going to be hidden, therefore not generating any interaction.

3.11 Failure checks

- **Issue** An OCE crashing or otherwise not performing correctly could break the running software.
- **Cause** Sometimes an implementation of a new feature will make the application crash. Failure checks then make sure to shut the experiment down when this happens. [19, 28]
- **Problem** Crashing software brings possible downtime and therefore possible loss of customers and revenue.

3.12 Number of changes

Issue The increasing complexity of experimenting with more than one change at once.

- **Cause** Usually an experiment consists of a version A and B. However sometimes an experimenter might want to experiment with multiple changes at the same time. For example, by giving every button on a page a specific color. However, these tests are usually more difficult and should be double checked if there is the possibility to see if the experiment can be done in the A/B fashion, [13] as this takes less time and effort.
- **Problem** This increased complexity increases the chance that statistical tests or something else in the experiment goes wrong

3.13 Day of week effect

Issue Not running an experiment for a limited amount of time.

- **Cause** The audience of software application can differ on time and day. On weekends different segments might visit the website [30, 53]. To represent the actual population of the website, it is important to run experiments for full weeks. This should even occur when the experiment becomes significant quickly, as running the experiment for the full week might result in the results not being significant.
- **Problem** If one section of the users is not correctly taken into account, this will give an incomplete result and possibly result in the wrong conclusion.

3.14 Not confirming the winning variant

Issue Not checking the version put into production.

- **Cause** After an experiment has concluded, the winning version is usually put in production and the experiment cleaned up.
- **Problem** When an experimenter does not make sure the winning version of the experiment works as expected, this could result in a broken version of the feature being put into production.

3.15 Encourage more experiments

- Issue Not asking for new experiment ideas at the right time.
- **Cause** An experiment has the goal to see if a hypothesis is true. However, while examining/discussing the results interesting insights or ideas for new experiments can be discovered
- **Problem** By not providing the space for an experimenter, the business might miss a great new idea.

3.16 Higher level question

- Issue Not being able to see the larger picture.
- **Cause** Experiments are not run in a vacuum. Most of the time they are used to answer a higher business level question [19].
- **Problem** If no infrastructure in place to think about these questions, it might be the case that the product ends up in a local optimum. It is therefore good to be able to group these experiments, so that this relation is clear to other users.

3.17 Share learnings

Issue Not being able to share learnings of experiments

Cause One factor important when becoming a data driven organization is to share the learnings from experiments. [19, 17, 35]. These can range from how to setup the best experiment to sharing the actual learnings from the experiment and having a repository of all experiments ever run [28]. From the interviews in chapter section 2.3, we also found that experimenters really feel a need for both.

Problem Experimenters will make the same mistake or miss out on great new product ideas

3.18 Rerun experiment when results are marginal

- **Issue** If the difference (significance/results) between two versions is low, rerunning the experiment could shed more light on which version performs better.
- **Cause** If the difference between two versions is small, i.e. if the metrics do change, but only by a little bit, then the experiment should be rerun again to validate that the results were correct and not the result of chance. [19]
- **Problem** Only barely significant results might be because of chance, not of actual difference.

3.19 Validation of experiment

Issue Not validating the experiment

- **Cause** Reproducibility is one the most important factors of science. If rerunning an experiment results in a different outcome, the setup was incorrect, it might be by chance or there are other influences which were not accounted for. Therefore, experiments which ran on the platform should be reproducible [19] and some experiments should be rerun on a certain timescale.
- **Problem** If experiments are not rerun to check the validity of experiments, this could lead to wrong experiments being performed.
Chapter 4

Solutions to pitfalls

In this chapter, we will analyze the problem and propose solutions to answer the second part of research question 2: "What are key pitfalls and their potential solutions?" (1.4).

4.1 Existing AB-test frameworks

To create solutions for the pitfalls it would be of help to look at other existing frameworks to see if any other platform has already created a solution to this problem and if so how they do this. We mainly look at other commercial products, as they have public documentation and online demos available. For most internal tooling, Microsoft's platform Exp [23] for Example, no screenshots or documentation could be found and therefore are not mentioned here.

4.1.1 Optimizely

Optimizely [53] is a popular AB-testing framework for performing experiments on websites. Optimizely claims that 24 of the top Fortune 100 companies are using Optimizely to perform tests. It uses Javascript to change the content of a page on the fly[42]. This enables it to be deployed quickly to any website, however Javascript has some downsides that should be discussed.

The first downside of using Javascript is that it takes time to load in the code to change it. As the content of the page usually renders earlier than the code that changes it, the user first sees the original page for a split second before it is changed. This is perceived as a flicker and therefore could influence the result of the experiment.

The second downside of using Javascript is that the code still needs to be loaded in, therefore making the page load slower. Third, this approach might not work on some browsers, as some features which this code needs might not be supported.

4.1.2 VWO

VWO [4] is another popular choice in the world of AB-testing, being used by Ubisoft and others. It uses the same Javascript-injecting method of creating variants. It is interesting to

note that none of the pitfalls we research are adressed in this tool.

4.1.3 Adobe Target

Adobe Target [8] is an Adobe Product focussed on online experimentation. Its documentation is extensive and clear about 4 of the pitfalls we are researching (Falsifiable hypothesis, Direction of change in hypothesis, minimum effect size and minimum duration). This gives the impression that the tool has been made with best practices in mind.

4.1.4 Google

From Google, we will look into 2 tools which enable experimenters to experiment: Firebase and Google Optimize.

Firebase

Firebase is a tool to "help mobile and web app teams succeed" [1]. It offers a lot of tools and services to help teams more easily and quickly create apps. One of the services they offer is A/B Testing. Firebase is the only website to offer a demo, which enabled to get a better understanding on how their features can solve our list of pitfalls.

Google Optimize

Google Optimize is part of Google's marketing platform to help engage customers "like never before" [2]. Although it is a separate product, many of the UX/UI elements seems very similar to Firebase.

4.2 Solutions

Given our overview the problems developers face and of the existing solutions, we can start proposing solutions to these problems. As some solutions solve multiple problems at once, each solution starts with a list of pitfalls it solves. After this, we use the existing solutions as inspiration to, finally, propose our solution to the pitfall.

Pitfall	Name of solution
Falsifiable hypothesis (01b, 3.1)	Enforce cross checks
Direction of change in hypothesis (02b,	Enforce cross checks
3.2)	
Guardrail metrics (02c, 3.3)	Add guardrail metrics
Technical debt (03a, 3.4)	Add checklist
Competitor safety (03b, 3.5)	Add checklist
Churning users (03c, 3.6)	Add checklist
Minimum effect size (04a, 3.7)	Enforce the correct experiment duration
Minimum duration (04b, 3.8)	Enforce the correct experiment duration

Pitfall	Name of solution
Withholding results (04c, 3.9)	Blur screen of the experiment
Simultaneous experiments (05a, 3.10)	Add checklist
Failure checks (06a, 3.11)	Add guardrail metrics
Number of changes (09a, 3.12)	Enforce cross checks
Day of week effect (11b, 3.13)	Enforce that experiments are run for full weeks
Confirm winning variant (12b, 3.14)	Add review step in experiment office
Encourage more experiments (13a, 3.15)	Add review step in experiment office
Higher level question (14b, 3.16)	Add review step in experiment office
Share learnings (15a, 3.17)	Add review step in experiment office
Rerun experiment when results are marginal (20a, 3.18)	Add button to rerun experiment
Validation of experiment (20b, 3.19)	Add button to rerun experiment

Table 4.1: Table showing the pitfalls and which solution is created to solve the pitfall. The first number behind the pitfall is the id of the pitfall, the second the section in which the pitfall is explained

4.2.1 Enforce cross checks

Resolves:

- Falsifiable hypothesis (1b)
- Direction of change in hypothesis (2b)
- Number of changes (9a)

Making sure the setup of the experiment is correct takes two steps, creating the setup and checking the validity of the setup. The Experiment Office has already taken great steps to help with the first part; By providing a template Figure 4.1 it is much easier to create a good hypothesis. However, hypotheses can still be made better and especially the direction of change is easy to forget, as "the number changing" completes the sentence, but is not valid. To solve this problem, it is proposed to have experimenters review eachothers hypothesis, similarly as is being done at Booking.com [28]. This review process can filter out many wrong or suboptimal hypotheses experimenters create.

Therefore we propose the following solution: We solve this problem by having other people check the experimental setup. Before your experiment can start, it is looked over by a group of experts, who could find mistakes in the experimental setup or tips on how to improve the experiment. The experiment is not allowed to start before the ok is given.

4.2.2 Add guardrail metrics

Resolves:

• Guardrail metrics (2c)

4. Solutions to pitfalls

ashboard Experimen	nt Overview	Experi	ments	Squad Ranking	Search	Documentation	ı		
ome > Experiments > Crea	ate experiment								
xperiment C	onfigur	ato	r						
1 Setup			2	The test		3 1	arget group		>
Title and key	Experiment t	itle	(j)						
	Experiment	(ey	1						
				Value is used by th	e developer:	s. Valid characters	are: "a-z0-!	9"	
	Feature		(j)			~	+ Cr	eate new feature	
Hypothesis ()	Template	Free-t	ext						
	Ву			replacing a drop-o	own with a	category page			
	will lead to			sales increasing					
	measured by			more orders					
Variants Control	Кеу	(j)	contr	ol-group					
	Description	(j)							
									8
Variant	Key	(j)	varia	nt1					
	Description	(j)							
									8
			+	Add variant					

Figure 4.1: Screenshot of the second step of the experiment creation workflow

• Failure checks (6a)

This solution is heavily based on the solution from Firebase (Google) in Figure 4.2, as it offers a very nice way of giving an overview of different metrics at the same time. Similarly, Airbnb [37] has a view (in Figure 4.3) to look at metrics which maybe less clear, but enables the experimenter to dive deeper into the data.

The view from Firebase should be tried first, as it has a lot of similarities to a study Meeuwsen et al. performed at ING [39]. Meeuwsen held a focus group with participants at ING. This resulted in the following sketch, seen in Figure 4.4, made by experimenters at ING on what a possible future version of the Experiment Office could look like.

This results in the following solution in the end:

To resolve this problem, we add the possibility for adding guardrail metrics and have the experimenter choose which metrics they also want to measure besides the default metric.

4.2.3 Add checklist

Resolves:

- Technical debt (3a)
- Competitor safety (3b)
- Churning users (3c)

Solutions to the above problems could not be found in existing solutions. One possible reason for this could be that these problems are not so much a constant problem, but do need attention once in the entire workflow. Therefore it seeems that a simple checklist could be sufficient to resolve these pitfalls.

Therefore we propose the following solution: To solve this, we show a quick list of things that you should think about when running experiments, we remind you of things that can go wrong while experimenting

4.2.4 Enforce the correct experiment duration

Resolves:

- Minimum effect size (4a)
- Calculate minimum duration (4b)

To solve this pitfall, data is needed, like the amount of visiters to a page and the size of the change they want to measure. Many online sample size calculators already exist which can help.

The solution from Optimizely is shown in Figure 4.5. It is able to determine the population size needed. The solution from Adobe Target in Figure 4.6 also calculates the number of days an experiment should be running. Therefore it seems this solution is best. Interestingly, both solutions are standalone and not integrated, something that could be a feature which the Experiment Office could integrate.

Continue runnin	inning g experiment to increas	e certainty		
Best perfor	ming variant: <u>Fev</u>	<u>w Ads</u> has a <u>51% chanc</u>	<u>e</u> of outperforming Base	Users in experir 2.1K
Details 」式 Re	emote Config 📄 Started	d Apr 24, 2017) 100% of users matching 1 criteria	🕼 4 variants
provement overview @)			
Variant	level_start 🕐 😭	user_engagement ②	Retention (4-7 days) ⑦	Crash-free users ③
Control group 561 users	Baseline	Baseline	Baseline	Baseline
More Ads 542 users	+0% -30% to +42%	+0% -100% to +167,024%	+0% -100% to +44,821%	+0.000% -0.196% to +0.192%
Less Ads 474 users	-10% -38% to +30%	+0% -100% to +172,989%	+8% -100% to +48,547%	+0.000% -0.220% to +0.190%
Few Ads	+1%	+0%	+0% -100% to +44.889%	+0.000%
529 users				
beriment results			😭 leve	I_start 👻 All variant
beriment results	ers triggering level s	tart event	♀ ieve	All variant
b29 users eriment results Percentage of use	ers triggering level_s	tart event	♀ leve	L_start → All variant s 30 days All time
Percentage of use	ers triggering level_s	tart event	♀ leve	I_start ▼ All variant s 30 days All time
Percentage of use 10% 5%	ers triggering level_s	tart event	♀ leve	I_start All variant a 30 days All time
Percentage of use 10% 5% 0% Jul 17	ers triggering level_st	tart event	♀ leve	L_start All variant a 30 days All time Jul 22
Percentage of use 10% 5% 0% r Jul 17 Variant	ers triggering level_s Jul ¹ 18 Improvement ©	tart event Jul 19 Jul 20 Probability to beat Probability to beat var	♀ leve	I_start ▼ All variant s 30 days All time Jul 22
529 users beriment results Percentage of use 10% 5% 0% 0% Jul 17 Variant Control group 561 users	Jul 18 Baseline	tart event Jul 19 Jul 20 Probability to beat Pro baseline 29	Q leve Q Q Ieve Jul 21 bability to be best Conversion (95%) 6	L_start → All variant s 30 days All time Jul 22 Irate level_start © 13.0% 58
529 users beriment results Percentage of user 10% 5% 0% 0% Jul 17 Variant Image: Control group 561 users Image: Set users Image: Set users Image: Set users	Jul 18 Improvement ③ Baseline +0% -30% to +42%	tart event Jul 19 Jul 20 Probability to beat Probability to beat Baseline 29' 50% 30'	Jul 21	L_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56
529 users beriment results Percentage of users 10% 5% 0% <td>Jul 18 Jul 18 Improvement ③ Baseline +0% -30% to +2% -10%</td> <td>tart event Jul 19 Jul 19 Jul 20 Probability to beat Probability and the set of the se</td> <td>Q leve</td> <td>I_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56 12.0% 44</td>	Jul 18 Jul 18 Improvement ③ Baseline +0% -30% to +2% -10%	tart event Jul 19 Jul 19 Jul 20 Probability to beat Probability and the set of the se	Q leve	I_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56 12.0% 44

Figure 4.2: Screenshot of the Google Firebase website. Screenshot taken from https://console.firebase.google.com/u/0/project/fir-demo-project/c onfig/experiment/results/3?hl=en

	Metric	Global	control			trea	tment			
	Filter metrics	Coverage	Mean	Mean	Percent	Change	p.	Value		MDE
~	Metric 1 Q Q targemetric	97.14%	0.05 430k/9.29M	0.05 429k/9.29M	0%		**** 0.5 [.]		-	0.73%
*	Metric 2 Q Q target metric		0.6 258k/430k	0.6 258k/429k	≠ 0.05% ≠ 0.38		**** 0.8i			2.4%
*	Metric 3 Q Q Core metric	96.54%	0.17 1.58M/9.29M	0.17 1.59M/9.29M			**** 0.4			1.2%
*	Metric 4 Q Q core metric	95.28%	0.06 577iu9.29M	0.06 580k/9.29M	0.41% ±1.4	\wedge	**** 0.5			2.1%
*	Metric 5 Q Q core metric		0.33 140k/430k	0.33 140k/429k	≠ 0.32% ± 0.64		**** 0.3			4.3%
*	Metric 6 Q Q core metric	96.33%	2.5 23M/9.29M	2.5 23M/9.29M			单杂录 1.0			1.1%
•	Metric 7 Q Q		0.32 566k/1.75M	0.32 565k/1.74M	0.62% ± 0.58		*** 0.04			1.9%
٠	Metric B Q @	42.47%	0.00 2.5769.29M	0.00 2.55k/9.29M		/	***	·~~		7.8%
	Metric 9 Q Q core metric See tickets , See issue breakdown	77.13%	0.00 23.4k/9.29M See tickets	0.00 23.4k/9.29M See tickets		\frown	victor 0.97			2.7%
*	Metric 10 Q @	85.01%	0.03 3196/9.29M	0.03 320k/9.29M	0%	\sim	*** 0.0	\sim	*	0.69%
*	Metric 11 Q @		4.7 1.56M/331k	4.7 1.54M/328k	▼ 0.05% ± 0.03	~	** < 0.0	\sim	*	0.24%

Figure 4.3: Screenshot of the metrics dashboard from Airbnb



Figure 4.4: Ideal dashboard, brainstormed by ING engineers during a focus-group performed by Meeuwsen et al.



Figure 4.5: Sample size calculator from Optimizely

The solution from Google, seen in Figure 4.7, is different, as they do not require an end date and just wait for the experiment to gain significance. However this generally considered a bad thing, as every experiment will become significant in the end, however the results will be negligible. Furthermore, from the interviews in section 2.3 we learned that currently only one experiment can be done on a page, so this solution is not chosen.

Therefore we propose the following solution: To aid people with setting the correct duration of the experiment, we introduce a screen where users are able to fill in all the information they have. How many people visit the page that is being experimented on, how detailed the result must be, etc. This will then automatically determine how many users need to participate in the experiment and therefore how long the experiment should last.

4.2.5 Blur screen of the experiment

Resolves:

• Witholding results (4c)

Adobe Target					
Sample Size Calcul	ator				
Conversion Rate Metric RPV Me	tric				
Confidence Level 95 %		Statistical Power 80 %		Baseline Conversion Rate	: (Control Offer)
Total Number of Daily Visitor	s	Number of Offers Including Con 5	trol	Daily Number of Visitors	per Offer
Lift (that can be detected wit	h power (80%) probability)		5%	10%	17.5 %
Absolute Difference in Conver	sion Rate (that can be detected w	vith power (80%) probability)	0.59%	1.18%	2.07%
Conversion Rate of Alternative	(that can be detected with power	er (80%) probability)	12.39%	12.98%	13.87%
Sample Size per Offer (# of vis	itors)		47,942	12,234	4,114
Sample Size per Offer (# of con	nversions)		5,657	1,444	485
Days to Complete Test			24	7	3
Weeks to Complete Test			4	1	1
Correct for Multiple Offe	rs (Bonferroni Correction)				

Figure 4.6: Sample size calculator from Adobe Target

To make sure experimenters do not become biased towards a solution, we can blur the results of the experiments. This solution is not found in any examples of other experiment platforms, as this (obviously) also hides the result of the experiment. To give the experimenter the ability to check if the experiment is running as expected, the blur can be removed.

If this solution is not chosen, a less drastic solution, which partially could solve the problem could be derived from Google's Firebase implementation (see Figure 4.2). By always displaying the chance that a variant performs better next to the improvement of the variant, it increases nuance and reduces the chance of people becoming biased to a variant. After all, saying that "a variant has a 100% better performance" gives a different feeling than "a variant has a 100% better performance, but there is a 52% it actually performs this well".

To conclude we propose the following solution: To resolve this problem, we can blur screen until the required number of user and significance has been reached. To aid users, we unblur the screen after users click yes on a pop-up where they are explained why looking at results would be a bad idea

Continue runnin	inning g experiment to increas	e certainty		
Best perfor	ming variant: <u>Fev</u>	<u>w Ads</u> has a <u>51% chanc</u>	<u>e</u> of outperforming Base	Users in experir 2.1K
Details 」式 Re	emote Config 📄 Started	d Apr 24, 2017) 100% of users matching 1 criteria	🕼 4 variants
provement overview @)			
Variant	level_start 🕐 😭	user_engagement ②	Retention (4-7 days) ⑦	Crash-free users ③
Control group 561 users	Baseline	Baseline	Baseline	Baseline
More Ads 542 users	+0% -30% to +42%	+0% -100% to +167,024%	+0% -100% to +44,821%	+0.000% -0.196% to +0.192%
Less Ads 474 users	-10% -38% to +30%	+0% -100% to +172,989%	+8% -100% to +48,547%	+0.000% -0.220% to +0.190%
Few Ads	+1%	+0%	+0% -100% to +44.889%	+0.000%
529 users				
beriment results			😭 leve	I_start 👻 All variant
beriment results	ers triggering level s	tart event	♀ ieve	All variant
b29 users eriment results Percentage of use	ers triggering level_s	tart event	♀ leve	L_start → All variant s 30 days All time
Percentage of use	ers triggering level_s	tart event	♀ leve	I_start ▼ All variant s 30 days All time
Percentage of use 10% 5%	ers triggering level_s	tart event	♀ leve	I_start ▼ All variant s 30 days All time
Percentage of use 10% 5% 0% Jul 17	ers triggering level_st	tart event	♀ leve	L_start All variant a 30 days All time Jul 22
Percentage of use 10% 5% 0% r Jul 17 Variant	ers triggering level_s Jul ¹ 18 Improvement ©	tart event Jul 19 Jul 20 Probability to beat Probability to beat var	♀ leve	I_start ▼ All variant s 30 days All time Jul 22
529 users beriment results Percentage of use 10% 5% 0% 0% Jul 17 Variant Control group 561 users	Jul 18 Baseline	tart event Jul 19 Jul 20 Probability to beat Pro baseline 29	Q leve T day: Jul 21 bability to be best Conversion (95%) 6 8.0%	L_start → All variant s 30 days All time Jul 22 Irate level_start © 13.0% 58
529 users beriment results Percentage of user 10% 5% 0% 0% Jul 17 Variant Image: Control group 561 users Image: Set users Image: Set users Image: Set users	Jul 18 Improvement ③ Baseline +0% -30% to +42%	tart event Jul 19 Jul 20 Probability to beat Probability to beat Baseline 29' 50% 30'	Jul 21	L_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56
529 users beriment results Percentage of users 10% 5% 0% <td>Jul 18 Jul 18 Improvement ③ Baseline +0% -30% to +2% -10%</td> <td>tart event Jul 19 Jul 19 Jul 20 Probability to beat Probability and the set of the se</td> <td>Q leve</td> <td>I_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56 12.0% 44</td>	Jul 18 Jul 18 Improvement ③ Baseline +0% -30% to +2% -10%	tart event Jul 19 Jul 19 Jul 20 Probability to beat Probability and the set of the se	Q leve	I_start → All variant s 30 days All time Jul 22 rate © level_start ⑦ 13.0% 58 13.0% 56 12.0% 44

Figure 4.7: Screenshot of the Google Firebase website. Screenshot taken from https://console.firebase.google.com/u/0/project/fir-demo-project/c onfig/experiment/results/3?hl=en

4.2.6 Enforce that experiments are run for full weeks

Resolves:

• Day of week effect (11b)

Both Optimizely [53] and Google Optimize [3] mention the "cyclical variations in web traffic"[3]. However, both do not seem to either force experiments to be full weeks or nudge experimenters to run their experiments for full weeks.

Therefore we propose our own solution: To resolve this problem, we present an experimenter with a warning when they want to stop an experiment on another weekday than the weekday the experiment was started. We add the options to let the experiment continue to the original end date, the next full week or stop it immediately.

4.2.7 Add review step in Experiment Office

Resolves:

- Confirm winning variant (12b)
- Encourage more experiments (13a)
- Higher level question (14b)
- Share learning (15a)

Optimizely has an extensive guide on how to convey the results of experiments to other teams mention the following basic items that should be in a report when conveying experiment results, copied directly from the webpage [41] :

- **Purpose**: Provide a brief description of "why" you're running this test, including your experiment hypothesis.
- **Details**: Include the number of variations, a brief description of the differences, the dates when the test was run, the total visitor count, and the visitor count by variation.
- **Results**: Be concrete. Provide the percentage lift or loss, compared to the original, conversion rates by variation, and the statistical significance or difference interval.
- Lessons Learned: This is your chance to share your interpretation of what the numbers mean, and key insights generated from the data. The most important part of results sharing is telling a story that influences the decisions your company makes and generating new questions for future testing.
- **Revenue Impact**: Whenever possible, quantify the value of a given percentage lift with year-over-year projected revenue impact.

This gives us a good overview of areas that need effort to keep track of the learnings of experiments (15a) and trigger experimenters to think of new experiments (13a).

To make sure the winning variant is working as expected, we can ask the experimenter to check the page after they cleanup the code to create the experiment. Inspiration to improve this conformation can be found with Optimizely on the page to preview variants in Figure 4.8. When creating experiments in Optimizely, there is an option to preview the



Figure 4.8: Preview of the experiment given by Optimizely

change inside the editor. Being able to preview a webpage in the tool can be very helpful and save the experimenter some time.

To make sure people keep in mind the higher level goals, they can be displayed in the list of metrics, like is proposed with guardrail metrics in subsection 4.2.2. This will help, but by explicitly asking experimenters about it we force them to think about it each time an experiment is wrapped up. It is most likely that when a the best variant to the experiment moves the team towards the goal there is not much to say here (although this should be further researched), but when a chosen version moves the team away from a goal, it might be useful to evaluate on this. (Is this the better solution? Is perhaps the goal wrong?).

Therefore we propose the following solution: After the experiment has concluded, we ask the experimenter to respond to a couple of questions related to the experiment they ran. What was their conclusion? What did they learn about experiments in general? What will be their followup experiments? How does this relate to their overall goal? Did they clean up their experiment and check if it was working as expected?

4.2.8 Add button to rerun experiment

Resolves:

- Rerun experiment when results are marginal (20a)
- Validation of experiments (20b)

No existing solutions for these problems exist in our survey of the tools. Therefore we propose our own solution: A button to rerun the experiment is added, which the experimenter can press if he/she would like to rerun their experiment to double check the results. This button is made more prominent if the significance of the experiment was low or if the experiment is sampled to check if the experiments are still running correctly

Chapter 5

Research methods used

In previous chapters we have determined pitfalls are a problem to experimentation and that 19 out of the 35 known pitfalls are relevant to this research. Next we seek to answer three questions. First, do experimenters indeed experience these pitfalls? Secondly, do they believe the new solution solves the problem? And lastly, are they hindered by the introduction of the solution? To answer these questions, we need to contact the experimenters. By creating a survey and sending this out to experimenters using the platform, we can find this information.

Surveys are chosen over interviews, as this results in a larger population to get more information from. Based on the preliminary interviews, a sufficient overview exists of the workflows at ING. If the results of the survey require follow-up, extra interviews/conversations can always be planned with specific individuals.

To perform the survey we will use Survalyzer[45], an internal ING tool, which respondents will be familiar with.

5.1 Survey contents

5.1.1 Introduction questions

As with any survey, a baseline is needed to understand the participant and get their permission to use the data from the survey for this master thesis.

We asked the participant for their email, in case the response raised questions we wanted to followup on, and the time the participant has been experimenting, to be able to notice any influence this has on the type of answers given.

5.1.2 Pitfall page

There exist 19 pitfalls we want to explore. For each of these pitfalls, a page exists within the survey. Every page contains 2 sections, where each section tries to answer a part of research question 2. Every question has been designed to not bias the participant and be clear to increase the reliability of measures [50]. Furthermore, we tried to make sure each question gives us as much information in as little time for the respondent. The last requirement

is achieved by having only 2 textboxes, one to collect information about the solution we propose and another as an optional catch all for any thoughts the participants might have, for example that they think the pitfall actually is not a problem. Both are optional, so that if the participant does not have enough time, they can skip it, as we would rather have that they give their opinion on more pitfalls, than explain one in great detail.

The first section of the page devoted to a pitfall explores the experimenter's view on the pitfall and tries to answer the first part of RQ2, namely how the experimenter experiences the pitfall. We want to know the prevalence of a pitfall, asking if they experienced it with a Yes/No question. Similarly, we ask about the prevalence of the knowledge of the pitfall. The perceived severity of the pitfall is gathered by using a Likert scale.

The second section details the proposed solution discussed in section 4.2 and tries to answer the second part of RQ2, how the experimenter experiences the solution. By asking participants if the solution solves the pitfall (Yes/No) and if it would help or hinder the participant with experimentation (again with in a Likert scale), we get an indication of the feasibility of the solution, which is then strengthened by the open text question of "Why?". This question is optional.

As discussed before, we wrap up with a catch all textbox, where we ask the participants if they have any other thoughts about the pitfall. If they do not have anything to add, they can skip this question.

Every page concerning pitfalls will look the same, except the description and solution of the pitfall. The contents of these can be found in Appendix C, together with screenshots of the survey.

5.1.3 Incorporating real-world examples

We would like to get as detailed responses as possible. If we can find out instances where an experimenter has fallen into a pitfall or encountered one, we can ask them detailed questions about this instance and their experience.

Only data for the day of week pitfall (3.13) can be extracted from the data of the experiment platform. This is accomplished by simply looking at the start and end date of an experiment and seeing if they are the same weekday. For the rest of the pitfalls, unfortunately, no instances can be extracted. This has a simple reasons. To solve pitfalls solutions need to be made, which do not exist yet. Only then will it become possible to gather which experiments have fallen for the pitfall. Secondly, metrics are currently not integrated in the platform, making tracking anything not already in the platform next to impossible.

Once collected, Survalyzer[45], the survey tool, has the ability to show participants links to only the experiments they performed, to make maximum impact.

5.2 Participation

We aimed for participants to complete at least 5 pitfalls. Each pitfall-part consists of 7 questions with the two initial questions. This totals to 35 questions. This number is large enough to gather information, but small enough that people are more likely to respond to

the survey, as it won't take forever. If, after this, participants want to give feedback on more pitfalls, we give them the opportunity to do so.

The ordering of the questions is as follows. If we have evidence that a participant has fallen for day-of-week pitfall, we start with that, with a personalized survey-section. Once they have completed it, we randomize the remaining pitfalls, as there is no pitfall we deem more important to get feedback on. This means that in the end, all pitfalls should have around the same number of participants.

Creating a list of participants turned out to be more difficult than expected. In the Experiment Office, the focus is put on the experiment and not the user who is executing the experiment. Because of this, there is no way to find out which user is executing a specific experiment (or any experiment). With help from the squad that maintains the Experiment Office, a list of active experimenters, who used the Experiment Office in the last 3 months, was extracted from log files. This list was further expanded with extra information. Each experimenter needs to be a part of a team to run experiments (2.3). These teams can be extracted. For each of these teams, if they created an experiment in the past 2 years, an effort was made to find the corresponding squad in the organization. To be a good citizen of ING, people in squads, for whom it was deemed very unlikely that they would have ever performed an experiment were excluded from participation. All the others were asked to participate in the survey. This resulted in 129 possible participants.

Chapter 6

Results

The survey ran from the 10th of June until the 22nd of June. In total 52 people responded (which is a response rate of 40.3%). Collectively, they gave feedback on 305 pitfalls. This chapter discusses the result of this survey. First in section 6.1, the overview of the results will be discussed, after which in section 6.3 (pitfalls) and section 6.4 (solutions) the results will be more deeply discussed and as a result, research question 2, "How are pitfalls and their solutions experienced?".

6.1 Overview of results

The distribution of answers, which can be seen in Table 6.1 seems to be evenly distributed. Due to the fact that the 'Full-week effects' pitfall was only shown to people who experienced this, this pitfall only received 5 responses. This is a result of both the response rate for the generic pitfalls being higher than expected and the response rate from the people, who experienced this pitfall being lower than expected. The experience with online controlled experimentation can be seen in Figure 6.1. On average participants had 3.67 years of experience.

Before diving deeper into the results, a few generic comments can already be made. For all the pitfalls, there are more people familiar with this pitfall than those who have actually experienced it, indicating that experimenters are proactively thinking about things which can go wrong when experimenting. There are six people, who said they experienced a pitfall, without being aware it existed. In two cases, this could have been as the participants did not think it was an actual problem. Others did not elaborate on this. For only the day of week effect this ratio does not hold, as we showed only people this question, for whom we were sure experienced this pitfall.

People on average think the new features are helpful, giving any feature an average of 3.52 (out of 5). This is because the entire population rates higher, with one outlier saying 9 solutions are 'very helpful' and only one being 'neither hindering, nor helpful'.

The following sections will look further into the results. For each of the questions a diagram was created showing the results to this question:

• Figure 6.2, which percentage has experienced a pitfall

6. RESULTS

Id	Name	Answers	Awareness	Experience	Solutions fixes problem
01b	Falsifiable hypothesis	17	58.82%	17.65%	58.82%
02b	Direction of change in hypothesis	17	58.82%	35.29%	64.71%
02c	Guardrail metrics	16	68.75%	25.00%	75.00%
03a	Technical debt	13	76.92%	76.92%	53.85%
03b	Competitor safety	16	56.25%	18.75%	50.00%
03c	Churning users	16	81.25%	43.75%	43.75%
04a	Minimum effect size	17	64.71%	35.29%	94.12%
04b	Calculate minimum duration	21	85.71%	52.38%	100.00%
04c	Witholding results	16	50.00%	25.00%	56.25%
05a	Simultaneous experiments	17	88.24%	41.18%	52.94%
06a	Failure checks	21	57.14%	23.81%	76.19%
09a	Number of changes	14	78.57%	57.14%	64.29%
11b	Day of week effect	5	60.00%	100.00%	60.00%
12b	Confirm winning variant	17	58.82%	17.65%	35.29%
13a	Encourage more experiments	18	77.78%	38.89%	61.11%
14b	Higher level question	14	57.14%	35.71%	50.00%
15a	Share learning	16	81.25%	68.75%	62.50%
20a	Rerun experiment when results are marginal	17	70.59%	47.06%	70.59%
20b	Validation of experiments	17	76.47%	29.41%	76.47%
	All	305.00	69.18%	38.69%	64.59%

Table 6.1: Table of all results







6.2. Differences among pitfalls and their solutions

Figure 6.2: Percentage of participants who experienced the pitfall

- Figure 6.3, which percentage is aware the pitfall exists
- Figure 6.4, a divergent stacked bar chart [12] of how severe the participant considers this pitfall
- Figure 6.5, percentage of participants who believes the solution solves the pitfall
- Figure 6.6, if the participant thinks the problem helps or hinders their work

These diagrams will be used to support the rest of this section.

Following the sample size calculation, with a confidence level of 95% and a 5% margin of error, 97 people should have responded. Therefore, the results of this survey should be seen as a starting point for further exploration, but cannot be used for drawing conclusions of statistical significance.

6.2 Differences among pitfalls and their solutions

The differences between the answers to pitfalls and solutions in the survey could be the result of random chance. To determine if this was the case or if the pitfalls are truly experienced differently or if some solutions are truly better than others, we employ statistical tests, which compare the answer from each pitfall/solution to all of the others.

The data collected with the survey is ordinal data [9]. This means that the our values are categories with an order. Therefore simple averaging is not possible.

6. Results







Figure 6.4: Divergent stacked bar chart [12] of the perceived severity of the pitfall



6.2. Differences among pitfalls and their solutions

Figure 6.5: Percentage of participants who believes the solution solves the pitfall



Figure 6.6: Divergent stacked bar chart [12] of the perceived usefullness of the solution

6.2.1 Cliffs delta

First, to determine if comparing two samples even makes any sense, the effect size is calculated using Cliff's delta [48]. Four different categories exist, Negligable (< 0.11), Small (0.11 - 0.28), Medium (0.28 - 0.43) and Large (> 0.43). We will only look at results having a large effect size (> 0.43).

6.2.2 Mann-Whitney U

The Mann-Whitney U test is a test developed by Mann, Henry B. and Whitney, Donald R. to determine the probability that two samples are significantly different[32]. For the calculation, the python library scipy was used [7].

We will be describing the pitfalls with a large effect size and where the significance of the Mann-Whitney U test is smaller than 0.05.

6.2.3 Results

All of the data used in the calculations for this section can be found in Appendix D.

For the question if the experimenter is aware of the problem, section D.1 shows that no comparison actually results in a large effect size. Therefore none of them have both a significant outcome and large effect size.

For the experience question this is different. As we only asked experimenters the "Fullweek effect" if they experienced it, this is significant and has a large effect size compared to other pitfalls (Table D.2). As we we only asked people who experienced it, we will not count it as a result. "Technical debt" is the second most significant result with 7 out of 18 pitfalls.

For severity, from Table D.2 it can be clearly seen that "Competitor Safety" is significantly not important to experimenters, as it shows a significant difference with 17 pitfalls.

For the question about solving the pitfalls, there is no real significant result. "Calculate minimum duration" has the best result and is significant with 7 other pitfalls (Table D.2).

For the HelpHinder question, there is another significant result found in Table D.10. Both "Minimum effect size" and "Calculate minimum duration" are significant with 13 other pitfalls. This is interesting for two reasons. First, this question revolves around the same solution, which makes it likely that they love this solution and secondly that this is a really great solution that will help experimenters.

6.3 Pitfall results highlights

This section will give some highlights on the awareness and perceived severity of the pitfalls. This section will answer the first partion of research question 2, namely how experimenters experience pitfalls.

• From Figure 6.2, it can be determined that "Technical debt", after "the Full weeks effects", is the most experienced pitfall. As many of participants are working with a CMS system, this is surprising. The open text does not explain this. Possible answers

could be that with a CMS technical debt still exists or that more coders than expected filled in the survey.

- It is interesting that "Falsifiable hypothesis", together with "Confirming winning variant", is the lowest experienced pitfall (Figure 6.2). Do participants think every one of their hypotheses has been falsifiable? Future work can investigate this further if this is the case.
- In Figure 6.3, even for the least known pitfall, 50% of the participants are aware of it. This indicates that knowledge about pitfalls is being shared or that they are researching themselves what can go wrong when experimenting.
- The perceived severity of "Competitor safety" is incredibly low (see Figure 6.4). Participants explain that either you should out-experiment your competitor, when running your experiment, or you are already production ready and that experiments can be used as a smoke screen to hide what you're actually working on.
- Unsurprisingly, customers leaving ING is the most perceived severe pitfall (see Figure 6.4)

6.4 **Results per solution**

This section explains what experimenters think of the proposed solutions and answers the second part of research question 2, how experimenters experience the solutions to the pit-falls.

6.4.1 Enforcing Cross checks

Solution: 4.2.1 - **Pitfalls:** Falsifiable hypothesis (3.1), Direction of change in hypothesis (3.2), Number of changes (3.12)

The solution seems to be very badly received. Although it is not the lowest scoring in the HelpHinder question ("Will this solution help or hinder you when running experiments?"), it is the pitfall where the most people said that the solution would hinder them. Most of the feedback in the "Why?" question focusses around the idea of having the experimenter being in control. Secondly they mention that they don't want anything delaying the running of their experiments. One participant even says there are already too many hurdles to go through. Another makes the point that experiments can be very unique to a business area and that even an expert in experimentation at ING might not understand the nuances of a particular experiment, which leads to delays in the running of an experiment.

6.4.2 Adding Guard rail metrics

Solution: 4.2.2 - Pitfalls: Guardrail metrics (3.3), Failure checks (3.11)

The solution seems to be well received. For both the Failure checks and Guardrail metrics pitfall, 76.2% and 75.0% think this solves the problem and most believe it will help. The textual answers are not very relevant in this case. In some cases it seems it was not clearly explained enough and the respondent is misunderstanding the solution. Other

answers agree that this is a good solution or commenting on who should choose which metrics to use; the team, defaults or a combination.

6.4.3 Adding a checklist

Solution: 4.2.3 - **Pitfalls:** Technical debt (3.4), Competitor safety (3.5), Churning users (3.6)

This particular solution is very interesting. The solution seems simple, however implementing seems to have problems. For all of the pitfalls this solution tries to solve, only 50% of the participants think it actually solves the problem. Furthermore, this solution tries to solve the most perceived severe pitfall (Churning Users) and least perceived severe pitfall (Competitor Safety) (Figure 6.4). All this results in the open-text answers to vary wildly. Some participants think that a checklist could prevent often made mistakes. Others believe that experimenters don't read, disagree on the usefulness of a generic checklist or think that they will remember anyways.

6.4.4 Enforcing the correct experiment duration

Solution: 4.2.4 - Pitfalls: Minimum effect size (3.7), Calculate minimum duration (3.8)

This solution was *very* well received. For the solution for the Minimum effect size almost all respondents perceived the solution to (slightly) help. Both pitfalls combined, only one participant (out of 38 who gave feedback) did not think this solution solved the pitfall. In the open text fields, many proclaim happiness that they no longer would have to use online calculators or spreadsheets to calculate effect size and experiment duration.

6.4.5 Blurring the screen of the experiment

Solution: 4.2.5 - Pitfalls: Witholding results (3.9)

The reactions to this solution are very mixed. Most do think that the solution solves the problem (56.25%), but how much it helps is contested. One participants writes down "Perfect solution", while others disagree, saying that this solution is annoying and unprofessional and that keeping an eye if the experiment is working correctly is necessary. One participant even goes so far to say that looking at progress gamifies the process, making him more enthusiastic to do more experiments.

6.4.6 Enforcing that experiments are run for full weeks

Solution: 4.2.6 - Pitfalls: Day of week effect (3.13)

This solution was not reviewed by as many participants as others, as discussed before. The response to this feature cannot easily be summarized. One unfortunate thing is that respondents do not seem to fully understand the problem and why the solution would/could work. One participant mentioning that they run their experiment for one week already, even though best practice is at least two weeks. Another mentioning something about changing the distribution afterwards. Another not understanding how stopping an experiment after two weeks and one day could influence the results. This could and should be better explained to experimenters. One experimenter does point out that even if you counteract this effect, there will still be others that influence the result of the experiments, like the time of the month when salaries are transferred.

6.4.7 Adding a review step in experiment office

Solution: 4.2.7 - **Pitfalls:** Not confirming the winning variant (3.14), Encourage more experiments (3.15), Higher level question (3.16), Share learning (3.17)

The reviews for this solution are overall pretty positive, although people think this solution is not very applicable to confirm the winning variant, where only 35% think this solves the problem, making this the worst solution in this area. The open text fields confirm this, as people do not seem to like the extra administration. Again, it seems that some participants did not fully recognize the problem.

6.4.8 Adding a button to rerun experiment

Solution: 4.2.8 - **Pitfalls:** Rerun experiment when results are marginal (3.18), Validation of experiments (3.19)

The response is positive, where 73.5% thinks the solution solves the problem. After "Enforce the correct experiment duration", participants think this solution will help them the most when performing experiments. Some respondents do want to nuance this to only rerun experiments when there are reasons to doubt the results, for example when the power is low. One respondent wants to use this to determine possible seasonality of changes. e.g. determine that variant A works better in the summer and B in the winter.

Chapter 7

Methodology for evaluating solutions in a real world scenario

There now are multiple solutions which can be implemented into the platform. Next answering research question 3, how one can evaluate a solution in a real-world scenario. Due to time constraints, this part could not be fully executed. However, work has been done and there are many lessons learned that are relevant, both for a future party wanting to build on this work and for ING. This chapter will elaborate on these lessons. This chapter will answer research question 3, how a solution can be evaluated in a real-world scenario.

7.1 Research methods used

In chapter 6, it was addressed which solutions experimenters prefer and which they would want to have implemented. This chapter addresses the question how such solutions cna be evaluated in practice. We aim to do this rigorously and in the spirit of continuous experimentation. This can take the form of a before/after study or a controlled experiment, where the difference is measured between the variants. This evaluation should be two pronged, which is discussed in the following sections.

7.1.1 Example

During the thesis project, work on the solution for the "Day of the week effect" (3.13) was already started. A screenshot of the solution can be found in Figure 7.1. When an experimenter fills in a date that is not a full week after the first date, it will pop up a warning to the user with more information and a button to fix it.

7.1.2 Making sure the solution works

When implementing one of the solutions, it is best practice to run an OCE to test if the solution works well. One of the most important parts is conducting the correct experiments.

In the case of the "Day of week effect", this would be how often an experimenter creates an experiment that does not run for full weeks, before and after treatment, but also how the

7. METHODOLOGY FOR EVALUATING SOLUTIONS IN A REAL WORLD SCENARIO

Dashboard	nt Overview Experiment	squad Ranking Sear	Documentation	
lome > Experiments > Cre	ate experiment			
xperiment C	onfigurator			
1 Setup		2 The test	3 Target group	\checkmark
Audience type	O Customers	s (ACMA id)		
	O Customers	s (RGB id)		
	O Employee	s (Corporate key)		
	Customers	s (Profile id / UUID)		
Start date	27-02-2020			
End date	11-03-2020			
O Your experiment d	oes not run for full weeks. <u>V</u>	Vhy is this bad? Fix it		×

Figure 7.1: Implementation of the full week effect solution

experiment resolves the pitfall. Do they click the button? Do they manually change it? Will they look for more information (click the link)? The full tracking plan can be seen in Appendix B. For each of these pitfalls one of these plans should be made and carefully executed.

Metrics were attempted to be implemented and lessons learned from this attempt are documented in section 7.2.

7.1.3 Making sure the solution is helping the experimenters

To determine if the solution has no unintended side effects, there should be a way for users to report these to improve the solution. Furthermore, if experimenters have improvements for the solutions, they should be able to explain this to the Tetris team (2.2). In section 7.3, lessons learned from using a survey to gather this type of data is explored and it is explained why this type of survey should not be used in this context.

7.2 WebTrekk

To be able to collect metrics, a tool is needed which is able to collect store and analyze this data. As ING is already using WebTrekk, it is the obvious choice to use. WebTrekk is an online analytics platform, the name of which will be changed to Mapp in the future [33]. WebTrekk is one of the main ways ING tracks the behavior of users accessing any of the online software ING offers.



Figure 7.2: Workflow of tracking a user on the platform

7.2.1 Obstacles with implementing WebTrekk

Implementing WebTrekk resulted in two main obstacles. The combination of the webplatform of ING and the infrastructure of ING resulted in the decision to reduce the scope of the thesis.

This highlights a general pain point. To become a data-driven organization, having tracking and metrics is a minimal requirement. To encourage the adoption measurement, the barrier to entry should be as low as possible and waiting for three months is the opposite of this. To encourage tracking, a team supporting and supervising this effort should be created.

Web Platform of ING

ING is currently in the process of moving from AngularJs [20], an old framework to build web applications, to Polymer [21], a framework no longer in development made by Google to use web components on the web.

The library for working with WebTrekk has been created for Polymer and WebComponents, which Angular is not compatible with. The experiment platform, however, is written in Angular. This made it particularly hard to integrate our work. There was an alpha version available, which should integrate with any Javascript based framework. However, the error messages were nondescriptive. This lead to finding the source of the error, the infrastructure of ING, too late for the possibility of finishing this software before the end of the internship.

The infrastructure of ING

The flow of tracking a user can be seen in Figure 7.2. A customer performs an action (clicking a button, loading a page), which has an event handler, which the WebTrekk library hooks into. A WebTrekk wrapper has been incorporated into many of the Polymer elements to make integration easier. As WebTrekk is third party software, the data needs to be sent to their servers. However it cannot be sent by ING directly for two reasons.

The first reason is that WebTrekk is a service which is not hosted on ING servers or at an ING url. Therefore calling this URL within a web page breaks Cross Origin Resource Sharing (CORS) [5] rules. It is a very good security feature and should not be tampared with. However, this results in the tracking library not being able to send tracking data to the WebTrekk servers.

Another reason not to send data directly to WebTrekk is to have a "circuit breaker" in place when an application breaks and sends a lot of (wrong) information to the WebTrekk servers. Instead of getting a large bill from WebTrekk, they can block the application and wait for it to recover or be fixed.

The solution to overcome this is to use of a reverse proxy. In this case, this results in the website sending data to the ING servers, which in turn sends it to the WebTrekk servers, as ING servers do not have to abide by the CORS rules.

This means that for every application at ING a reverse proxy has to be set up. This requires a change in the configuration of the network of ING, which for security purposes needs to go through a certain workflow. Estimates for the length of this procedure turned out to be around 3 months and this made it impossible to implement in the timeframe left.

7.3 The platform survey

To determine if there are any problems with the proposed solution, a feedback point needs to be created. We could ask experimenters if they are satisfied with the new solutions, however this could have unintended consequences. Asking if experimenters are satisfied with the new solutions will always involve some prejudice. The experimenter might give positive feedback to not disappoint the person asking for it or come up with feedback on the spot, which is not the type of feedback we want. After all, if the experimenter does not care (enough), we resolve the pitfalls without any change in user satisfaction. Furthermore, with the survey we already have a good indication if the solution will be positively received. To counter these possible problems in the feedback cycle, we created a generic survey, which is described below. Because it is a generic survey, the Tetris team (2.2) can determine if there are other aspects of the tool which can be improved and therefore it is more likely that the team allows such a feature to be introduced in the platform.

7.3.1 The survey layout

The survey has two requirements. The first requirement is to determine if the overall satisfaction has not gone down. This is satisfied by the first question on the survey, which is a Likert scale asking how satisfied the developer is with the experiment platform. This question gives us an opportunity to look at the satisfaction over time and therefore give us a reasonable metric of satisfaction over time. The term reasonably here is used as opinions change over time and by definition are subjective. However, if there are enough results, we should be able to average these subjective data points out and get a somewhat objective result.

The second requirement is to determine if the given satisfaction given by the experimenter in the first question has anything to do with the solutions which have been implemented. To do this we present the experimenter with a textbox where they can write detailed feedback on what they like and don't like about the platform. To try and get the experimenter to talk about the solutions the following placeholder text is used: "Write as

Feedback form How satisfied are you with the experiment platform?						
O Very dissatisfied	o	о ОК	o	O Very satisfied		
Please write as m	orate on yo	our choice?				
We love your inpu						

INC Foodbook Form

Figure 7.3: Mockup of the survey

much as possible" and "We love your input". The intent is that this will increase the length of the feedback. If the experimenter does not mention the new solutions, this will be either because they deemed it not important enough to talk about, they are happy with it or they do not care. If they do mention the new solutions, we should check if their feedback is valid and actionable. If that is the case, we could improve the solution further. Otherwise we note down this result. If the feedback is positive, we can also note this. Lastly, we ask the experimenter if we can have their email address. This might be of value if we want to follow-up with the experimenter to get a better understanding of why they do not like a solution or if the feedback they gave is not clear. This is not a required field because of course this survey is anonymous to make sure that we do not bias the results by having the experimenter think there might be some negative consequences if they criticize the platform. A mockup of the survey is shown in Figure 7.3 and the final survey is shown in Figure 7.4.

7.3.2 Placing the platform survey

There are three places in the workflow where this survey was placed. The first one is after the experiment-creation workflow and can be seen in Figure 7.5. The creation workflow is when the experimenter has created the experiment in the experiment platform. When creating a new experiment, the experimenter must go through a few steps to set the name, hypothesis and target group. At the end of this workflow we ask experimenters what they think of the platform. This has another benefit, as many solutions are placed in the creation workflow. Therefore, this is a great place to ask for feedback.

The second place a survey will be created is after the experiment has concluded, see

7. METHODOLOGY FOR EVALUATING SOLUTIONS IN A REAL WORLD SCENARIO

How satisfied are you with the experiment platform? Please select the appropriate smiley
Could you please elaborate on your choice?
Please write as much as possible. We love your input ;-) 0/1000
Can we have your email address to contact if we want to followup with you?
Optional
Submit

Figure 7.4: Final survey

ING M The Guide Dashboard Experiment Over Home > Experiments > Create expe	Fruitloops A/B Testing view Experiments Squad Ranking Search Documentation riment
xperiment Conf	igurator
Example implementation	<pre> Use the following code to make frontend variants:</pre>
Finish	hat do you think of the experiment platform? We would love to know your feedback! Take the survey >

Figure 7.5: Button to the survey

Figure 7.6. At this point the experimenter has possibly checked the website during the runtime of the experiment, looking at the intermediate result, but possibly forgot that experiment was running. To remind the experimenter that the experiment has concluded, we, in collaboration with the Tetris team, created a feature that will send an email to the developer telling them that their experiment has concluded. This opportunity is then also used to ask the experimenter to fill in a survey. Again, this is a great place to ask for feedback, as they have gone through the entire flow of an experiment and therefore can tell exactly where the pain points of the system are and what can be improved.

Lastly, there is a generic feedback button. On every page, as can be seen in Figure 7.5, a user can press the "Give your feedback" button to give feedback on the platform.

Asking at these locations in the process does have some drawbacks. The survey after the experiment-creation-workflow leaves out the demographic of people, which left the workflow halfway through. As we ask experimenters about their opinion after the workflow has finished, this does not include the people who do not finish the workflow (obviously). Unfortunately, it is impossible and impractical to ask a person to fill in a survey when leaving the website and therefore we cannot reach this demographic. The same holds for the survey shown in the reminder email. We might have lost experimenters, who have stopped their

7. METHODOLOGY FOR EVALUATING SOLUTIONS IN A REAL WORLD SCENARIO



Figure 7.6: Reminder email

experiments midway through.

7.3.3 Preliminary Results

As only few people responded to the survey, we will briefly discuss the results below.

On the 14th of February the first respondent filled in the survey. At the time of writing, July, only 9 experimenters have filled in the survey. 5 of these are just praising the experiment platform, usually saying "great" and rate the system good (2) or great (3). 3 respondents fill in the survey asking for the feature of tracking more than one metric and give a fair (2) or bad (1) rating. The last respondent wanted technical help and gave the platform a bad rating.

7.3.4 Conclusion on the use of a survey

Although setting up this survey was promising, the results are preliminary and should not be used to validate a possible implementation in the end, as they will not have enough results to make definitive conclusions on the solution. These results do indicate that people want more metrics, something that can and should weigh in when deciding the roadmap of the platform.

Chapter 8

Conclusion & Discussion

8.1 Conclusion

Online controlled experimentation is a powerful tool, but can be deceivingly easy, as experimenters can face many pitfalls while performing their experiment. It is vital for experimenters and organisations to perform correct and valuable experiments. To aid them, this thesis tries to answer the question: "*How can we reduce human error in online controlled experimentation*?" To answer this question, exploratory interviews were held to determine common problems experimenters face. By looking at common pitfalls developers face, based on previous work, a list of possible solutions was created, taking existing solution as starting points. A survey, answered by 52 experimenters at ING, was performed to validate the result of the interviews and to evaluate the proposed solutions on the perceived helpfullness to the experimenter.

8.1.1 RQ1 What are key pitfalls and their potential solutions?

By looking at existing literature, we established that a list of 19 pitfalls are relevant for this thesis. By leveraging existing solutions, such as Optimizely [53], VWO [4], Adobe Target [8] and Firebase [1] as inspiration, solutions can be created to solve these problems. This thesis proposes 8 solutions to solve the 19 pitfalls, ranging from invasive changes, such as adding a review step to the workflow, to simple solutions like a checklist.

8.1.2 RQ2 How are pitfalls and their solutions experienced?

Experimenters are well aware of problems that can occur while experimenting. During interviews, we determined that experimenters are already using statistical methods to determine the minimum runtime of experiments. The survey confirmed that for each pitfall more than 50% of the experimenters is aware that this problem could present itself. Less experimenters have actually experienced the pitfalls. Experimenters think each pitfall to be roughly equally severe with the exception of "Competitor safety", which they consider not severe.

There are mixed reactions to the proposed solutions. The solutions were viewed in two different perspectives: whether it solves the problem and whether the solution helps the experimenter. Experimenters are immensely enthusiastic about the solution "Enforcing the correct experiment duration", solving "Calculate the minimum duration" and "Calculate effect size". Respondents believed the solution solved these problems 100% and 94.1% respectively and almost all respondents perceived the solution to (slightly) help the experimenter in performing their experiments. The next best solution (in solving the problem) drops by 20%, making other solutions still relevant to look at, but not as much. The same holds for the perceived usefullness of the solution, where some ("Rerunning experiments" and "Validation of experiments") are highly rated, but not as much. There are not truly badly received results. Worst received is the proposed solution solving "The number of changes" pitfall, where more than half of the respondents thought the solution (slightly) hindered the experimenter.

RQ3 How can one evaluate a solution in a real-world scenario?

From the exploratory work done in this thesis, we determined that evaluating a solution in a real-world scenario runs into technical challenges. To evaluate a solution it must be possible to track and log the interaction an experimenter has with the platform, and to receive feedback from experimenters. Both are not adequately implemented into the platform. Tracking on the platform does not exist and implementing tracking runs into technical and bureaucratic problems. Current feedback methods do exist and work, although there is not enough feedback to be able to draw conclusions. Therefore, at this moment, there is no way of evaluating a solution in a real-world scenario.

Summary

There are a lot of pitfalls experimenters face while performing online controlled experiments. This thesis shows that experimenters are well informed about the existence of pitfalls and believe that almost all should be resolved. There are many promising solutions to these pitfalls which experimenters rate as helpful. However evaluating these solutions at this point in the current context is not possible.

8.2 Implications

8.2.1 Implications for Industry

This work is an example for how an organization can start with implementing continuous experimentation. The case study at ING shows what experiments are being performed and can inspire readers for new avenues of experiments. The main work provides an overview of which pitfalls experimenters deem important and which solutions will help experimenters. This provides any company an indication which pitfalls are most important to resolve first and which solutions could help them in resolving these pitfalls.
8.2.2 Implications for ING

Based on the results of the survey, there now is a list of possible improvements which can help experimenters with experimentation. It should then be no surprise that one of the recommendations is to start implementing the solutions.

In chapter 7, a lot of shortcomings of the process are described and what hurdles still remain in the data driven culture of ING. These should be resolved, so that any employee can improve the interactions ING has with their customers.

Finally, during the time of the internship and research a view on how the experiment office can be improved further was created, which is discussed in Appendix A.

8.2.3 Implications for Research

To the authors' knowledge, this is the first body of knowledge surveying experimenters with solutions for pitfalls, before implementation. This shows interesting insights into which pitfalls are important. Besides points mentioned in 8.2.1, this thesis also tries to provide a jumping off point for any new research being performed under the banner of the AI for Fintech Research (AFR).

8.3 Threats to validity

8.3.1 Internal validity

Internal conditions can affect the independent variable with respect to causality, without the researcher being aware [50]. In regard to this thesis this includes any influences on the survey and interview, which we cannot know.

One possible threat is the threat of experimental mortality, i.e. the loss of participants during the experiment. The survey was sent out to many people. People with limited experience with continuous experimentation might not be motivated to start or continue with the survey if they believe they cannot give relevant answers and/or are not passionate about experimentation. Experienced experimenters might be more motivated to fill in the survey, resulting in biased results.

Exploration into this threat resulted Figure 8.1, which shows the distribution of answers not to be clearly biased towards a particular experience group.

8.3.2 Construct validity

Construct validity concerns generalizing the result of the experiment to the concept or theory behind an experiment [50].

For the survey, we presented the participants with solutions to pitfalls. However, in the eyes of participants, this might look like a possible feature list. The question if the solution would help or hinder might be also be seen as asking how much they want this feature. Answers to these question might then be biased. This bias might then also increase in areas where the experimenter has more experience. This threat cannot be mitigated, but can be taken and has into account when analyzing the results.



Figure 8.1: Diagram showing the combined number of responses for pitfalls, set to the years of experience

8.3.3 Conclusion validity

Conclusion validity threats can affect the ability to draw correct conclusion [50].

For this research, the biggest threat is low statistical power, as the number of people running experiments at ING is high enough that performing interviews is an insane amount of work, but performing a survey, even with the great response rate seen in this case, is not a large enough sample size to significantly determine the results. However, anything short of making the survey mandatory will not change this.

8.3.4 External validity

External valididity threats are conditions that limit the ability to generalize the results of the experiment industrial practice [50].

As the interviews and survey were performed at a single company, this introduces the possibility that this work is not generalizable over the industry, as views on pitfalls might differ, based on what the needs of the organization are. Furter research is needed to validate these results at other companies.

8.4 Future work

In the overall picture of continuous experimentation research at ING and the AI4Fintech research group, this research can serve as a foundation to a lot of other work.

8.4.1 Evaluate solutions in a real-world scenario

In chapter 7, it was concluded that there were significant hurdles to evaluating the features. These hurdles can be overcome with time. A future researcher can perform this work or wait until a time engineers at ING have performed this research.

8.4.2 Reproduce this research in different industries

This work has shown a first indication of the perceived severity and value of solutions to practitioners. To be truly able to say that experimenters value these solutions, this survey should be repeated at different companies, favourably in different (non FinTech) industries.

8.5 Ethics of online controlled experimentation

Online controlled experiments can bring enormous advantages to companies [53, 17, 30]. It is therefore no surprise that many [46, 15, 38, 36, 24, 16] use them. However, there are some serious ethics questions related to using experiments to try and sway user actions.

OCEs reveal an interesting fact about us, humans. Where one would love to pretend that we act with only logic, example after example in the world of online controlled experimentation shows us that this idea that we have come to believe is wrong. Colours can influence our click behaviour [52] and the change of an image increases the amount of mortgages bought [40].

In 2012, Facebook published a study showing that if a Facebook user was shown more positive posts, their own posts also showed more positive words. The same held for negative words. Another study showed that Facebook could increase voter turnout by displaying badges on users pages [51]. These results show both that there is enormous positive potential in the use of experimentation to increase happyness, but also gives us an insight into the enormous power these companies hold.

Furthermore, it is reasonable to assume that few of the participants of the experiments were aware they were being targeted (until 2012, it was not event mentioned in the terms of service of Facebook that your data could be used for research purposes [11, 51]).

In an academic experimental setting this raises serious questions [54]. While performing research, I could find no company performing experiments on their users actually disclosing this to them. In fact, I was surprised at the breadth of the companies performing OCEs and the scale at which these experiments were performed. It is hard to opt-out of something you do not know is happening. The new cookie law by the European Union seems to have had a positive influence in this case, as it is best practice to use cookies to track users [16] and

decide variants based on a user cookie. Therefore, by not accepting cookies, you also do not give consent being experimented on.

Another large problem in the area of online controlled experimentation seems that there does not seem to be an ethical code of conduct regarding experimentation [11]. There should guidelines to determine areas where OCEs are not ethical or even harmful to the customer of the product. For example, forbidding to perform experiments to increase the amount of children buying V-bucks with their parent's credit card, to tempt broke gamblers to to just play another round of poker or to increase sales of mortgages if this has long term negative effects for the customer.

As with any tool, they can be used for good and bad. As such, OCEs are no exception and therefore time and effort should be put in to make sure they are used correctly and ethical.

Bibliography

- [1] Firebase. URL https://firebase.google.com/. Accessed: 2020-05-18 02:03:24.
- [2] Website, a/b testing & optimization tools google optimize, . URL https://market ingplatform.google.com/about/optimize/. Accessed: 2020-05-18 02:11:58.
- [3] Vwo, . URL https://vwo.com/. Accessed: 2020-05-18 02:03:24.
- [4] Vwo. URL https://vwo.com/. Accessed: 2020-05-18 02:03:24.
- [5] Cross-origin resource sharing, 2017. URL https://www.w3.org/TR/2014/REC-c ors-20140116/.
- [6] Research contributions in human-computer interaction acm interactions, 2020. URL https://interactions.acm.org/archive/view/may-june-2016/resear ch-contribution-in-human-computer-interaction.
- [7] scipy.stats.mannwhitneyu scipy v1.5.2 reference guide, 2020. URL https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats .mannwhitneyu.html.
- [8] adobe.com. Personalization adobe inc. URL https://www.adobe.com/marketin g/target.html.
- [9] Alan Agresti. Categorical data analysis. Wiley-Interscience, Hoboken, N.J, 2013. ISBN 978-0-470-46363-5.
- [10] Florian Auer and Michael Felderer. Current state of research on continuous experimentation: A systematic mapping study. In 2018 44th Euromicro Conference on Software Engineering and Advanced Applications (SEAA). IEEE, aug 2018. doi: 10.1109/seaa.2018.00062. URL https://doi.org/10.1109%2Fseaa.2018.00062.
- [11] Raquel Benbunan-Fich. The ethics of online research with unsuspecting users: From a/b testing to c/d experimentation. *Research Ethics*, 13(3-4):200-218, nov 2016. doi: 10.1177/1747016116680664. URL https://doi.org/10.1177%2F 1747016116680664.

- [12] Mark Bounthavong. Communicating data effectively with data visualization part 15 (diverging stacked bar chart for likert scales) — mark bounthavong, 2020. URL https://mbounthavong.com/blog/2019/5/16/communicating-data-effecti vely-with-data-visualization-part-15-divergent-stacked-bar-chart-f or-likert-scales.
- [13] Thomas Crook, Brian Frasca, Ron Kohavi, and Roger Longbotham. Seven pitfalls to avoid when running controlled experiments on the web. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining -KDD '09.* ACM Press, 2009. doi: 10.1145/1557019.1557139. URL https://doi.or g/10.1145%2F1557019.1557139.
- [14] Pavel Dmitriev, Somit Gupta, Dong Woo Kim, and Garnet Vaz. A dirty dozen. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, aug 2017. doi: 10.1145/3097983.3098024. URL https://doi.org/10.1145%2F3097983.3098024.
- [15] eng.uber.com. experimentation-platform from eng.uber.com, 2020. URL https:// eng.uber.com/experimentation-platform/. (Accessed on Mon Sep 16 2019).
- [16] exp platform.com. Documents 2019-first practical online controlled experiments summit sigkdd explorations.pdf from exp-platform.com, 2020. URL https://exp-platform.com/Documents/2019-FirstPracticalOnlineContr olledExperimentsSummit_SIGKDDExplorations.pdf. (Accessed on Fri Sep 13 2019).
- [17] Aleksander Fabijan, Pavel Dmitriev, Colin McFarland, Lukas Vermeer, Helena Holmström Olsson, and Jan Bosch. Experimentation growth: Evolving trustworthy a/b testing capabilities in online software companies. *Journal of Software: Evolution and Process*, 30(12):e2113, nov 2018. doi: 10.1002/smr.2113. URL https: //doi.org/10.1002%2Fsmr.2113.
- [18] Aleksander Fabijan, Pavel Dmitriev, Helena Holmstrom Olsson, and Jan Bosch. Online controlled experimentation at scale: An empirical survey on the current state of a/b testing. In 2018 44th Euromicro Conference on Software Engineering and Advanced Applications (SEAA). IEEE, aug 2018. doi: 10.1109/seaa.2018.00021. URL https://doi.org/10.1109%2Fseaa.2018.00021.
- [19] Aleksander Fabijan, Pavel Dmitriev, Helena Holmstrom Olsson, Jan Bosch, Lukas Vermeer, and Dylan Lewis. Three key checklists and remedies for trustworthy analysis of online controlled experiments at scale. In 2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP). IEEE, may 2019. doi: 10.1109/icse-seip.2019.00009. URL https: //doi.org/10.1109%2Ficse-seip.2019.00009.
- [20] Google. Angularjs superheroic javascript mvw framework, URL https://angu larjs.org/.

- [21] Google. Polymer project, . URL https://www.polymer-project.org/.
- [22] Jayant Gupchup, Yasaman Hosseinkashi, Pavel Dmitriev, Daniel Schneider, Ross Cutler, Andrei Jefremov, and Martin Ellis. Trustworthy experimentation under telemetry loss. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, oct 2018. doi: 10.1145/3269206.3271747. URL https://doi.org/10.1145%2F3269206.3271747.
- [23] Somit Gupta, Lucy Ulanova, Sumit Bhardwaj, Pavel Dmitriev, Paul Raff, and Aleksander Fabijan. The anatomy of a large-scale experimentation platform. In 2018 IEEE International Conference on Software Architecture (ICSA). IEEE, apr 2018. doi: 10.1109/icsa.2018.00009. URL https://doi.org/10.1109%2Ficsa.2018.00009.
- [24] hbr.org. 2017 09 the-surprising-power-of-online-experiments from hbr.org, 2020. URL https://hbr.org/2017/09/the-surprising-power-of-online-experim ents. (Accessed on Fri Sep 27 2019).
- [25] Hennie Huijgens, Davide Spadini, Dick Stevens, Niels Visser, and Arie van Deursen. Software analytics in continuous delivery. In *Proceedings of the 12th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. ACM, oct 2018. doi: 10.1145/3239235.3240505. URL https://doi.org/10.1145%2F 3239235.3240505.
- [26] ING.com. 2019 annual report, URL https://www.ing.com/About-us/Annual-r eporting-suite/Annual-Report/2019-Annual-Report.htm.
- [27] ING.com. 2018 annual report, URL https://www.ing.com/About-us/Profile/ ING-at-a-glance.htm.
- [28] Raphael Kaufman, Jegar Pitchforth, and Lukas Vermeer. Democratizing online controlled experiments at booking.com. 10 2017.
- [29] Ron Kohavi, Randal M. Henne, and Dan Sommerfield. Practical guide to controlled experiments on the web. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining KDD '07*. ACM Press, 2007. doi: 10.1145/1281192.1281295. URL https://doi.org/10.1145%2F1281192. 1281295.
- [30] Ron Kohavi, Roger Longbotham, Dan Sommerfield, and Randal M. Henne. Controlled experiments on the web: survey and practical guide. *Data Mining and Knowledge Discovery*, 18(1):140–181, jul 2008. doi: 10.1007/s10618-008-0114-1. URL https://doi.org/10.1007%2Fs10618-008-0114-1.
- [31] Ron Kohavi, Diane Tang, Ya Xu, Lars G. Hemkens, and John P. A. Ioannidis. Online randomized controlled experiments at scale: lessons and extensions to medicine. *Trials*, 21(1), feb 2020. doi: 10.1186/s13063-020-4084-y. URL https://doi.org/10. 1186%2Fs13063-020-4084-y.

- [32] H. B. Mann and D. R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. Ann. Math. Statist., 18(1):50–60, 03 1947. doi: 10.1214/aoms/1177730491. URL https://doi.org/10.1214/aoms/1177730491.
- [33] Mapp. Webtrekk is now a full-fledged part of mapp. URL https://www.webtrekk .com/.
- [34] David Issa Mattos, Pavel Dmitriev, Aleksander Fabijan, Jan Bosch, and Helena Holmström Olsson. An activity and metric model for online controlled experiments. In *Product-Focused Software Process Improvement*, pages 182–198. Springer International Publishing, 2018. doi: 10.1007/978-3-030-03673-7_14. URL https://doi. org/10.1007%2F978-3-030-03673-7_14.
- [35] medium.com. Duolingo: 1 URL https://medium.com/googleplaydev/duo lingo-1-improvement-every-week-ab7d61689119. (Accessed on 14-Jun-20 12:09:38).
- [36] medium.com. airbnb-engineering experiments-at-airbnb-e2db3abf39e7 from medium.com, 2020. URL https://medium.com/airbnb-engineering/experim ents-at-airbnb-e2db3abf39e7. (Accessed on Mon Sep 16 2019).
- [37] medium.com. airbnb-engineering https-medium-com-jonathan-parks-scaling-erf-23fd17c91166 from medium.com, 2020. URL https://medium.com/airbnb-eng ineering/https-medium-com-jonathan-parks-scaling-erf-23fd17c91166. (Accessed on Mon Sep 16 2019).
- [38] medium.com. netflix-techblog selecting-the-best-artwork-for-videos-through-a-btesting-f6155c4595f6 from medium.com, 2020. URL https://medium.com/netfl ix-techblog/selecting-the-best-artwork-for-videos-through-a-b-tes ting-f6155c4595f6. (Accessed on Mon Sep 16 2019).
- [39] Juliette Meeuwsen, Matthijs van Leeuwen, Hennie Huijgens, Kevin Anderson, and Arie van Deursen. Stimulating the adoption of a/b testing in a large-scale agile environment. 5 2019.
- [40] Ernst Mulder. Data driven decisions: Validating and supporting a continuous experimentation development environment. Master's thesis, TU Delft Electrical Engineering, Mathematics and Computer Science, 12 2019. https://repository.tudelft.nl/islandora/object/uuid%3A08f2c0b4-2a a8-4e12-9b58-073dcdfb4553?collection=education.
- [41] optimizely.com. Share your results with stakeholders, URL https://help.optim izely.com/Analyze_Results/Share_your_results_with_stakeholders.
- [42] optimizely.com. Optimizely x web experimentation, . URL https://www.optimize ly.com/platform/experimentation/.

- [43] Radiolab. The trust engineers—radiolab—wnyc studios, 2020. URL https:// www.wnycstudios.org/podcasts/radiolab/articles/trust-engineers. (Accessed on Tue Mar 24 2020).
- [44] Murali Krishna Ramanathan, Lazaro Clapp, Rajkishore Barik, and Manu Sridharan. Piranha: Reducing feature flag debt at uber.
- [45] Survalyzer. Professional survey tool and panel management survalyzer, 2020. URL https://www.survalyzer.com/.
- [46] Diane Tang, Ashish Agarwal, Deirdre O'Brien, and Mike Meyer. Overlapping experiment infrastructure. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '10.* ACM Press, 2010. doi: 10. 1145/1835804.1835810. URL https://doi.org/10.1145%2F1835804.1835810.
- [47] Edith Tom, Aybüke Aurum, and Richard Vidgen. An exploration of technical debt. Journal of Systems and Software, 86(6):1498 – 1516, 2013. ISSN 0164-1212. doi: https://doi.org/10.1016/j.jss.2012.12.052. URL http://www.sciencedirect.com/ science/article/pii/S0164121213000022.
- [48] András Vargha and Harold D. Delaney. A critique and improvement of the cl common language effect size statistics of mcgraw and wong. *Journal of Educational and Behavioral Statistics*, 25(2):101–132, 2000. doi: 10.3102/10769986025002101. URL https://doi.org/10.3102/10769986025002101.
- [49] Wikipedia. Covid-19 disease wikipedia. URL https://en.wikipedia.org/wik i/Coronavirus_disease_2019.
- [50] Claes Wohlin. *Experimentation in software engineering*. Springer, Berlin New York, 2012. ISBN 978-3-642-29044-2.
- [51] www.nytimes.com. A bright side to facebook's experiments on its users, 2020. URL https://www.nytimes.com/2014/07/03/technology/personaltech/the-bri ght-side-of-facebooks-social-experiments-on-users.html. (Accessed on 14-Jun-20 12:09:38).
- [52] www.nytimes.com. 2009 03 01 business 01marissa.html from www.nytimes.com, 2020. URL https://www.nytimes.com/2009/03/01/business/01marissa.htm l?pagewanted=print. (Accessed on Fri Sep 27 2019).
- [53] www.optimizely.com. customers from www.optimizely.com, 2020. URL https:// www.optimizely.com/customers/. (Accessed on Fri Feb 28 2020).
- [54] Sezin Yaman, Fabian Fagerholm, Myriam Munezero, Hanna Maenpaa, and Tomi Mannisto. Notifying and involving users in experimentation: Ethical perceptions of software practitioners. In 2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM). IEEE, nov 2017. doi: 10.1109/esem.2017.31. URL https://doi.org/10.1109%2Fesem.2017.31.

Appendix A

Future version of the experiment office

During the months working and thinking about improving the platform, insights was gained into what a future version of the platform could look like. The following is a personal opinion on what this future version could look like and has not been validated.

In Figure A.1, one can find this new version. The overall design is based on an accordion view, inspired by the Firebase UI [1]. When a user opens the page for an experiment, the page opens the step of the process the experimenter is at. If the experimenter wants to look at another step (for example Experiment Setup instead of the Results), they can click on that section and it will expand, to show all the information the experimenter has filled in, while the current section is hidden. This design tries to minimize the use of dialogs, such that any new experimenter can look at a random experimenter, see all the settings used, and be able to rerun it without much trouble.

This version also has better integration with WebTrekk, showing a dashboard, inspired by what Meeuwsen et al. found[39] (see Figure 4.4) and how the dashboard of Airbnb looks Figure A.2. This will be either a dashboard created in-house or the WebTrekk dashboard.

It is very important to have one single place to keep track of experiments. From literature, we have determined that it can be very useful to have a searchable repository of experiments [28]. This approach should be two fold. First, ING should be making an effort to consolidate experimentation to a single platform. Secondly, the platform should integrate with as many products as possible. By importing data, possibly read-only, the platform becomes the go to place if a person wants to know something about a (type of) experiment.

This could lead to the following workflow. John wants to see what effect changing the color of a button on the homepage has on the click-through-rate. John goes to the experiment office to look for anyone who has done experiments with buttons and colors in the last year and finds that his colleague, Hank, has done so on a different platform than John, but it shows up, because it is integrated. After a quick conversation, John learns that the color green usually works best. John creates his new experiment on the platform and creates the variants. Using the platform he determines how long the experiment should run for (6.4.4) and starts the experiment.

All experiments | All experiments | Global leaderboard Experiment Multi Does adding green to page x increase profils? Ends in 2 days Experiment Setup External Connection Results St.p Advice: Run your experiment loyer Shocass View Onta e WT Conclusion Teardown



	Metric	Global	control			treatment			
	Filter metrics	Coverage	Mean	Mean	Percent	Change	P-1	falue	MDE
*	Metric 1 Q Q target metric	97.14%	0.05 430k/9.29M	0.05 429k/9.29M	0%	\sim	भ्रेश्रम् 0.51		• 0.73%
•	Metric 2 Q @ target metric		0.6 258k/430k	0.6 258k/429k	▼ 0.05% ≠0.38		治治治 0.80		2.4%
*	Metric 3 Q @ core metric	96.54%	0.17 1.58M/9.29M	0.17 1.59M/9.29M			**** 0.42		/ 1.2%
*	Metric 4 Q @	95.28%	0.06 577/09.29M	0.06 580k/9.29M		<u> </u>	**** 0.58	~~~~	2.1%
*	Metric 5 Q @		0.33 140k/430k	0.33 140k/429k			**** 0.32	~~~~	4.3%
*	Metric 6 Q @	96.33%	2.5 23M/0.29M	2.5 23M/9.29M			前立放 1.00		/ 1.1%
*	Metric 7 Q Q		0.32 566k/1.75M	0.32 565k/1.74M	0.62% ± 0.58		*** 0.04	\frown	/ 1.9%
*	Metric B Q @	42.47%	0.00 2.57k/9.29M	0.00 2.55k/9.29M		\nearrow	**** 0.77	~~	7.8%
	Metric 9 Q Q core metric See tickets , See issue breakdown	77.13%	0.00 23.46/9.29M See tickets	0.00 23.4k/9.29M See tickets			*** 0.97		2.7%
*	Metric 10 Q @	85.01%	0.03 3196/9.29M	0.03 3206/9.29M	0%	$\sim \sim$	*** 0.08		• 0.69%
*	Metric 11 Q @ (1 day behind		4.7 1.56M/331k	4.7 1.54M/328k	▼ 0.05% ± 0.03		*** < 0.01	\sim	• 0.24%

Figure A.2: Screenshot of the metrics dashboard from Airbnb

Appendix B

B. TRACKING PLAN

Tracking plan





74





An experimenter opens the experiment office to leave feedback using the button at the top

ld				Description	Туре	Tag
Home		O Martjn (M.	n (MJW) Steenbergen Give your feedback Logout	User logs into the experiment platform	Page	expo:home
	ING M The Guide Fruitloops	8 Testing	f	They click on the feedback button	Event	expo:home.genericFeedbackButton
	Dashboard Experiment Overview Experim	ents Squad Ranking Search Documentation	totion Create an experiment >			
	> Home					
	Dashboard					
	Welcome Martijn (M.J.W.) Steenbergen! You are currently not associated with any squad. To To be added to a squad, please contact the Tetris Te	mmended to join a squad.				
	rollout	a/b test	a/b test			

Appendix C

Survey Information

This appendix gives an extensive overview of the survey. In Figure C.1, a screenshot is shown of the first page of the survey. In Figure C.2, one of the survey pages is shown, where questions are asked about a specific pitfall. For each of the pitfalls, everything is the same, except for two items. First, in the text of the first question, an explanation of the pitfall is. Then, secondly, halfway around the page, the explanation for the solution is put. The exact descriptions for each of these pages can be found in this appendix.

C.1 Pitfall descriptions

C.1.1 Falsifiable hypothesis (01b)

Explanation

A falsifiable hypothesis means that a hypothesis can be proven wrong. If a hypothesis cannot be proven wrong, the experiment has less value, as the result is already known. "The views on this page will increase" is a falsifiable hypothesis "The people will be happier with a green button" is not, because there is no way of measuring happiness with metrics from a web pag and therefore impossible to (dis)prove.

Solution

We solve this problem by having other people check the experimental setup. Before your experiment can start, it is looked over by a group of experts, who could find mistakes in the experimental setup or tips on how to improve the experiment. The experiment is not allowed to start before the ok is given.

C.1.2 Direction of change in hypothesis (02b)

Explanation

When creating a hypothesis for an experiment, usually it contains a description of what will happen to the metric that is measured. For example "When removing the sign up, we

C. SURVEY INFORMATION

i nank you for participating:
The layout of this survey is centered around problems you, as an experimenter, might face while experimenting. Each of these problems is its own page in the survey. We start with an explanation of the problem and ask you about your experience with the problem. Next we propose a solution to this problem and ask your expertise on whether this solution is fitting.
It would be great if you complete at least three pages . This should take you around 20 minutes. After these three sections you can continue answering more sections if you want. You can stop at any time . The progress will be saved automatically. If you want, you can even continue at a later time.
Before we start, we have a few questions for you.
 Conditions and Stipulations 1. I understand that all information is confidential and that I will not be personally identified, unless I agree to it. I agree to complete the survey for research purposes and that the data derived from this anonymous survey may be published in journals, conferences, and blog posts. 2. I understand that my participation in this research survey is totally voluntary and that declining to participate involves no penalty or loss of benefits. I may withdraw my participation at any time. I also understand that I may decline to answer any question shout the research answering. 3. I understand that I can contact the researchers if I have any questions about the research. I am aver that my consent may not directly benefit me. I am also aware that the author will maintain the data collected in perpetuity and may utilize data for future academic work. 4. By checking the "Confirm" box below I freely provide consent and acknowledge my rights as a voluntary research participant as outlined above and provide consent to the researchers to use my information in conducting research. I also confirm that I am of legal age to fill in this survey.
O Lagree
Can we contact you for followup questions? If yes, write down your email below)
How long have you been experimenting?
Next
It would be great if you could share your experience of five problems you faced. After these five problems you can continue answering more sections if you want. You can stop at any time. The progress will be saved automatically. If you want, you can eve continue at later time.

Figure C.1: The first page of the survey

79

Witholding randil Career constructions within the expected means the state of the means are balance of the state o					
Intermediate results are visible for experimenters. It right is not the the backening blace down of a solution. If writing the solution is not the	Witholding resul	lts			
like in the scale between the scale below in a base regarding the Witholding results problem (If intermediate resu	Its are visible for experimenters, it r	night lead to them becoming biased	l towards a solution. If version A st	arts out as better
	Before reading this,	were you aware this problem existe	:d?		
	O Yes	O No			
How severe do you think this problem is? Cotice and each to globe the handle In this hannonging This hannonging This hannonging This hannonging Our proposed solution to solve Withholding results Consciout this problem, we can build so screen with lither required number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of states and significance has been reached. To aid users, we unblue the second number of second number of states and significance has been reached. To aid users, we unblue the second number of second numb	Have you ever encod	untered this problem yourself?			
Not a problem This is annoying This right factor bases are bases are based on a possible of the problem of the p	How severe do you the scale to place	think this problem is?			
Our proposed solution to solve Witholding results To resolve this problem, we can blur screen until the required number of user and significance has been reached. To aid users, we unblur the cance of the series cick we on a pop-up where they are explained with looking at results would be a bad idea. Deveat thick this solution solves this problem? I've meet the place tableade the bad with teal badde. Provide the solution solves this problem? I've meet the place tableade the badde the badde teal badde. Provide the solution place tableade the badde teal badde. I've meet the place tableade the badde teal badde. Provide the solution place tableade the badde teal badde. I've meet the place tableade the badde teal badde teal badde. Provide the solution place tableade teal badde teal badde. I've meet to place tableade tableade teal badde teal badde. Provide the solution place tableade teal badde	Not a problem	This is annoying	This might cause issues when experimenting	This will cause issues when experimenting	Should be resolved as soon as possible
Our proposed solution to solve Witholding results To resolve this problem, we can blur screen until the required number of user and significance has been reached. To aid users, we unblur the screen after users click yes on a pop-up where they are explained why looking at results would be a bad idea Do you think this solution solves this problem? Ure any other they here leave? For any other they here leave? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the Witholding results problem? For any other thoughts you want to share regarding the for any other thought the problem? For any other thought the problem and					
Dry our think this solution solves this problem? I've model I've model How much would this feature help or hinder to do your experiments? Citize on the scale to place handle Hinder Slightly hinder Nor hinder nor help Help Are there any other thoughts you want to share regarding the Witholding results problem?	Our proposed so To resolve this probl screen after users cl	lution to solve Witholding res lem, we can blur screen until the rec lick yes on a pop-up where they are	ults juired number of user and significar explained why looking at results wo	nce has been reached. To aid users uild be a bad idea	, we unblur the
How much would this feature help or hinder to do your experiments? Click on the scale to place handle Hinder Slightly hinder Nor hinder nor Slightly help Hinder Slightly hinder Verbre? Control Contro Control <td>Do you think this so (If you answer no, please el Yes</td> <td>lation solves this problem? laborate in the why text below)</td> <td></td> <td></td> <td></td>	Do you think this so (If you answer no, please el Yes	lation solves this problem? laborate in the why text below)			
Hinder Slightly hinder Nor hinder nor help Slightly help Help Why? (Optional) Coptional Coptional Coptional Coptional Slightly help Help Help Help Slightly help Help Help Slightly help <td>How much would th Click on the scale to place</td> <td>iis feature help or hinder to do your : handle</td> <td>experiments?</td> <td></td> <td></td>	How much would th Click on the scale to place	iis feature help or hinder to do your : handle	experiments?		
Why? (Optional) Image: Contrast of the state in the state regarding the Witholding results problem? (Optional)	Hinder	Slightly hinder	Nor hinder nor help	Slightly help	Help
Why? (Optional) Image: Contract of the system of					
Are there any other thoughts you want to share regarding the Witholding results problem?	Why? (Optional)				
Are there any other thoughts you want to share regarding the Witholding results problem?					
Are there any other thoughts you want to share regarding the Witholding results problem?					
	Are there any other	thoughts you want to share regardi	ng the Witholding results problem?	,	
Back	Back				Nevt

It would be great if you could share your experience of five problems you faced. After these five problems you can continue answering more sections if you want. You can stop at any time. The progress will be saved automatically. If you want, you can even

expect signups to go down". In this case the metric measured is the number of signups and the direction is that this number will decrease. Writing this down forces the experimenter to explicitly hypothesize what will happen. Also statistically, the maths between calculating significance and calculating the possibility of a significant change in one direction is different than calculating the possibility of a significant change in one of two directions.

Solution

We solve this problem by having other people check the experimental setup. Before your experiment can start, it is looked over by a group of experts, who could find mistakes in the experimental setup or tips on how to improve the experiment. The experiment is not allowed to start before the ok is given.

C.1.3 Guardrail metrics (02c)

Explanation

Guardrail metrics enable the experimenter to keep track of extra metrics to make sure the experiments do not harm the overall goal of the business. An example of this is keeping track of the number of sales when running an experiment in which changes the color of a single button. If guardrail metrics are not in place, this could lead to unwanted results, even though the experiment is a success. To continue the example, the color might lead to more people clicking the button, but the sales going down.

Solution

To resolve this problem, we add the possibility for adding guardrail metrics and have the experimenter choose which metrics they also want to measure besides the default metric.

C.1.4 Technical debt (03a)

Explanation

Technical debt is the implied cost of additional rework in the future by choosing an easy (limited) solution now instead of taking a better approach that would take longer.

Solution

To solve this, we show a quick list of things that you should think about when running experiments, we remind you of things that can go wrong while experimenting

C.1.5 Competitor safety (03b)

Explanation

Running AB tests in production might give a competitor an indication of a new product (line) being developed and give competitors more time to catch up or beat time to market.

Solution

To solve this, we show a quick list of things that you should think about when running experiments, we remind you of things that can go wrong while experimenting

C.1.6 Churning users (03c)

Explanation

If, for example, an AB test gives customers a bad onboarding experience, they might not choose to join after all, making ING lose a customer.

Solution

To solve this, we show a quick list of things that you should think about when running experiments, we remind you of things that can go wrong while experimenting

C.1.7 Minimum effect size (04a)

Explanation

The minimum effect size enables you to, with statistics, determine how many users you need to test with to measure a 5%, 10% or 100% change. After all, it is easier to see if a new version is 100% better than 1%, as the cause of a 1% change is more likely randomness if you have a small number of participants. The minimum effect size is part of the calculation of how long the experiment should run. If this is not calculated, the experiment could run longer than necessary or too short to get actual results.

Solution

To aid people with setting the correct duration of the experiment, we introduce a screen where users are able to fill in all the information they have. How many people visit the page that is being experimented on, how detailed the result must be, etc. This will then automatically determine how many users need to participate in the experiment and therefore how long the experiment should last.

C.1.8 Minimum duration (04b)

Explanation

The minimum duration is how long an experiment should at least run to gain significance. If this is not calculated, the experiment could run longer than neccesary or too short too short to get actual results.

Solution

To aid people with setting the correct duration of the experiment, we introduce a screen where users are able to fill in all the information they have. How many people visit the

C. SURVEY INFORMATION

page that is being experimented on, how detailed the result must be, etc. This will then automatically determine how many users need to participate in the experiment and therefore how long the experiment should last.

C.1.9 Withholding results (04c)

Explanation

If intermediate results are visible for experimenters, it might lead to them becoming biased towards a solution. If version A starts out as better this might lead to "rooting" for this version to win.

Solution

To resolve this problem, we can blur screen until the required number of user and significance has been reached. To aid users, we unblur the screen after users click yes on a pop-up where they are explained why looking at results would be a bad idea

C.1.10 Simultaneous experiments (05a)

Explanation

If users are being exposed to multiple AB-tests at once, this could lead to invalid results. The most extreme example is when one test tests a different color for a button, while another test removes the button.

Solution

To solve this, we show a quick list of things that you should think about when running experiments, we remind you of things that can go wrong while experimenting

C.1.11 Failure checks (06a)

Explanation

Sometimes an implementation of a new feature will make the application crash. Failure checks then make sure to shut the experiment down when this happens.

Solution

To resolve this problem, we add the possibility for adding guardrail metrics and have the experimenter choose which metrics they also want to measure besides the default metric.

C.1.12 Number of changes (09a)

Explanation

Usually an experiment consists of a version A and B. However sometimes an experimenter might want to measure multiple things at the same time. For example by giving every button on a page a specific color. Test 1: Button A is green, Button B is yellow Test 2: Button A is red, Button B is also red

However, these tests are usually more difficult and should be double checked if there is the possibility to see if the experiment can be done in the A/B fashion, as this takes less time and effort.

Solution

We solve this problem by having other people check the experimental setup. Before your experiment can start, it is looked over by a group of experts, who could find mistakes in the experimental setup or tips on how to improve the experiment. The experiment is not allowed to start before the ok is given.

C.1.13 Day of week effect (11b)

Explanation

Every day a different segment of the users log on to the website. For example, during the week, people who work might not have time to look at a new mortage, while during the weekend they do. Starting and ending the experiment on a different weekday might change the distribution of the visitors of the page so much it is no longer representing the actual distribution of the users visiting the website.

Solution

To resolve this problem, we present an experimenter with a warning when they want to stop an experiment on another weekday than the weekday the experiment was started. We add the options to let the experiment continue to the original end date, the next full week or stop it immediately.

C.1.14 Confirm winning variant (12b)

Explanation

After an experiment has concluded, the winning version is usually put in production and the experiment cleaned up. When an experimenter doesn't make sure the winning version of the experiment works as expected, this could introduce unintentional bugs or changes.

Solution

After the experiment has concluded, we ask the experimenter to respond to a couple of questions related to the experiment they ran. What was their conclusion? What did they

C. SURVEY INFORMATION

learn about experiments in general? What will be their followup experiments? How does this relate to their overall goal? Did they clean up their experiment and check if it was working as expected?

C.1.15 Encourage more experiments (13a)

Explanation

An experiment has the goal to see if a hypothesis is true. However, while examining/discussing the results interesting insights or ideas for new experiments can be discovered. By not providing the space for an experimenter, the business might miss a great new idea.

Solution

After the experiment has concluded, we ask the experimenter to respond to a couple of questions related to the experiment they ran. What was their conclusion? What did they learn about experiments in general? What will be their followup experiments? How does this relate to their overall goal? Did they clean up their experiment and check if it was working as expected?

C.1.16 Higher level question (14b)

Explanation

Experiments are not run in a vacuum. Most of the time they try to answer a higher business level question. If no infrastructure in place to think about these questions, it might be the case that the product ends up in a local optimum. It is therefore good to be able to group these experiments, so that this relation is clear to other users.

Solution

After the experiment has concluded, we ask the experimenter to respond to a couple of questions related to the experiment they ran. What was their conclusion? What did they learn about experiments in general? What will be their followup experiments? How does this relate to their overall goal? Did they clean up their experiment and check if it was working as expected?

C.1.17 Share learnings (15a)

Explanation

Sharing experiences when experimenting can be very valuable to other experimenters. Think of helping experimenters answer the following questions: What are common problems? How can experimenters test better? If these are not shared, many can make the same mistakes, costing valuable time and effort from the experimenters.

Solution

After the experiment has concluded, we ask the experimenter to respond to a couple of questions related to the experiment they ran. What was their conclusion? What did they learn about experiments in general? What will be their followup experiments? How does this relate to their overall goal? Did they clean up their experiment and check if it was working as expected?

C.1.18 Rerun experiment when results are marginal (20a)

Explanation

If the difference (significance/results) between two versions is low, rerunning the experiment could shed more light on which version performs better.

Solution

A button to rerun the experiment is added, which the experimenter can press if he/she would like to rerun their experiment to double check the results. This button is made more prominent if the significance of the experiment was low or if the experiment is sampled to check if the experiments are still running correctly

C.1.19 Validation of experiment (20b)

Explanation

Reproducibility is one the most important factors of science. If rerunning an experiment results in a different outcome, the setup was incorrect or there are other influences which were not accounted for. Therefore experiments which ran on the platform should be reproducible.

Solution

A button to rerun the experiment is added, which the experimenter can press if he/she would like to rerun their experiment to double check the results. This button is made more prominent if the significance of the experiment was low or if the experiment is sampled to check if the experiments are still running correctly

Appendix D

Statistic tables

The following tables are an n-n comparison of the pitfalls.

D.1 Effect Size

For this work, for each combination of pitfalls the effect size was calculated [48]. See the shorter version in section 6.2. The values in bold have a large effect size (> 0.43).

205 0176471 0176471 00176471 00176471 00176471 001217730 001217750 01121777 011217777 011217777 011217777 011217777 011217779 011217779 011217779 011217779 011217779 011217779 010122777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 010127777 0101277777 0101277777 010127777777 010127777777777		20b 0.117647 0.1117647 0.0058824 0.0058824 0.106618 0.1045113 0.105618 0.105618 0.115647 0.1176471 0.1176471 0.1176471 0.1176471 0.1176471 0.1176471 0.1176471 0.1176471		20b 0.0734622 0.0739585 0.02739585 0.0272794 0.07272455 0.07272451 0.07272455 0.072795431 0.0716431 0.0716437 0.07176437 0.071779 0.071770 0.071770 0.071770 0.071700 0.071700 0.071700 0.07170000000000
20a 0117647 0117647 0017647 0017647 00053348 0106518 0106518 01076471 01076471 01076471 0107882 0107882 0107882 0107882 0107882 0107882 01076618 01006618 01006618		20a 0.294118 0.2294543 0.022058543 0.022058543 0.0220585 0.0230585 0.0230585 0.01250585 0.01250585 0.01250585 0.01250585 0.01250585 0.01250585 0.0125471 0.0115447 0.01154471 0.01154471		20u 0.128028 0.128028 0.1280286 0.2090586 0.2380529 0.2383235 0.235333 0.235333 0.23534 0.23534 0.225822 0.2238724 0.023882 0.2238724 0.0238754 0.0238754 0.0238754 0.0238754 0.0238754 0.0238754 0.0238754
154 0.234265 0.234265 0.1234265 0.0443269 0.0446441 0.0446441 0.0446441 0.044677 0.024755 0.024755 0.0106618 0.047794		15a 0511029 0511731 05341559 0437559 043755 0534599 0534599 01469405 014000000000000000000000000000000000		154 0.0025 0.01025 0.01025 0.0215677 0.0215677 0.0215677 0.0275059 0.0125765 0.0125765 0.0125765 0.0125765 0.0125765 0.0125765 0.012656 0.0125765 0.012656 0.027666 0.02766 0.0000000000000000000000000000000000
145 0.016807 0.016807 0.016807 0.016807 0.0198929 0.0310929 0.0310924 0.0028171 0.0028571 0.0028571 0.0028571 0.0028571 0.0028477		14b 14b 0.189672 0.0041203 0.00412038 0.1019143 0.10196453 0.101966457 0.101966457 0.101966457 0.0081357 0.0081357 0.0081358 0.01194628 0.01194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.001194638 0.0011948 0.00011948 0.000100000000000000000000000000000000		14b 0.063025 0.0085025 0.008659 0.1096599 0.1096599 0.1233214 0.1233214 0.1035211 0.0037231 0.0037231 0.0036612 0.1036902 0.1036902 0.1036502 0.12305382 0.01205382 0.02305382 0.033053 0.033055 0.033055 0.033055 0.03505 0.035055 0.035055 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.035555 0.0355555 0.0355555 0.035555555555
13a 0.189542 0.000278 0.000278 0.0008278 0.0084779 0.0084779 0.0187778 0.0187778 0.0187778 0.017977 0.017977 0.017977 0.017977 0.017977 0.017977 0.017977 0.017977 0.017977 0.017977 0.017972 0.0179772 0.0179772 0.0179772 0.0179772 0.0179772 0.0179772 0.0179772 0.0179772 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.0179777 0.01797777 0.01797777 0.01797777 0.01797777 0.01797777 0.01797777 0.01797777 0.01797777 0.01797777 0.017977777 0.0179777 0.01797777 0.01797777 0.01797777 0.017977777 0.017977777 0.01797777777 0.01797777777777777777777777777777777777		13a 0.212418 0.212418 0.212418 0.0035948 0.0035611 0.0153681 0.0153681 0.0153681 0.0153681 0.015378 0.015378 0.015378 0.015378 0.015374 0.023661 0.022661 0.0000000000000000000000000000000000		13a 0.156863 0.156863 0.156863 0.2186859 0.238589 0.2417967 0.2417967 0.2312859 0.23028869 0.1386869 0.1386869 0.118687 0.1865080808080000000000000000000000000000
2b 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	61	12b 0.176471 0.075529 0.075676 0.016471 0.075539 0.016471 0.0755339 0.061652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.001652 0.0016552 0.0016520 0.0016520 0.0016520000000000000000000000000000000000		12b 0.031142 0.031142 0.0283088 0.283408 0.283408 0.283408 0.28340 0.118507 0.118507 0.118507 0.118507 0.118504 0.018508 0.0106685 0.010685 0.010685 0.010685 0.000685 0
11b 0011765 0011765 0011765 0011765 002125 002225 002275 0022571 0022871 00285771 00285771 00285771 00285771 00285771 00285771 00285771 00285771 00285771 00285771 00285771 00277755 00277757 00277757 00277757 00277757 00277757 00277757 00277757 00277757 002777757 00277777 00277777 00277777 00277777 00277777777	xisted	11b 825529 0.557059 0.757059 0.757059 0.57019 0.57019 0.57019 0.588235 0.588235 0.588235 0.76619 0.76619 0.7588235 0.75619 0.611111 0.642857 0.53882 0.54882 0.55882 0	elf?	11b 0011765 0.0011765 0.001538 0.001538 0.01538 0.01538 0.0114256 0.01145666 0.01145666 0.0114566666666666666666666666666666666666
09a 0.1974.99 0.1974.99 0.0982.14 0.0982.14 0.0982.14 0.0164.84 0.0164.84 0.0174.89 0.0174.89 0.0174.85 0.0179.57 0.0179.57 0.01079.57 0.00757.86 0.01079.57 0.00079.57 0.000790000000000000000000000000000000	blem e	03. 0.304958 0.2184567 0.2184567 0.2184567 0.2184563 0.2184563 0.2384555 0.23845555 0.238455555555 0.2384555555555555555555555555555555555555	yours	0%a 0.0756.3 0.0756.3 0.0756.47 0.0758.71 0.0758.71 0.0758.75 0.058.44 0.0758.75 0.0258.75 0.028.83 0.014.286 0.014.
0068 0016807 0016807 01166071 01166071 0016807 00241071 00241071 00241071 00241071 00241071 00241071 00241071 00241454 00241071 0134454 00241071 0134454	iis pro	000 0001625 0.001625 0.0011905 0.001905 0.00190500000000000000000000000000000000	oblem	06a 0.002409 0.0022409 0.0022493 0.0022493 0.002493 0.0032493 0.0032493 0.0032493 0.0032493 0.0032493 0.003249 0.01325 0.013252 0.012252 0.002524 0.012252 0.012252 0.012252 0.002524 0.002545 0.002550 0.002500 0.00250000000000
0.254118 0.254118 0.254118 0.1544138 0.1544138 0.1544138 0.194835 0.198535 0.219853 0.219853 0.21924 0.282353 0.294418 0.282353 0.282353 0.282353 0.282353 0.282353 0.282353 0.282353 0.282353 0.294418 0.282353 0.282353 0.282353 0.294418 0.282353 0.294418 0.294453 0.294418 0.294553 0.294453 0.294453 0.294555 0.2945555 0.2945555 0.2945555 0.2945555 0.294555555555555555555555555555555555555	vare th	05a 0.2353294 0.2353294 0.2353294 0.005882-4 0.0161765 0.2357466 0.2357466 0.2357466 0.0225735 0.117647 0.0225735 0.0225824 0.02258245 0.02258245 0.02258245 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.0225824 0.025874 0.025874 0.025874 0.025874 0.025874 0.025874 0.022787476 0.022787470 0.022787470 0.022787470 0.0227874 0.0227874 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.022787470 0.0227874700000000000000000000000000000000	this pr	054 0.190311 0.190311 0.072874 0.0772876 0.072863 0.715863 0.070584 0.010284 0.010284 0.010284 0.0102842 0.212885 0.01005 0.212885 0.212885 0.01005 0.212885 0.01005 0.0000000000
046 0088235 0188235 0188235 0188235 0188235 018735 0147059 0147059 0138235 0138235 0138235 0138235 0138235 013871429 01387142000000000000000000000000000000000000	you av	04c 0.073529 0.073529 0.073529 0.073521 0.062530 0.062530 0.0187531 0.0187531 0.011965 0.011965 0.01138899 0.013889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.0138889 0.013889 0.013889 0.013889 0.013889 0.013889 0.0138889 0.013889 0.0138889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.0148889 0.014889 0.014889 0.014889 0.00489000000000000000000000000000000000	itered i	04c 0.128676 0.128676 0.3128676 0.32341 0.3331731 0.43375 0.4375000000000000000000000000000000000000
04b 0258908 0268908 0108643 0108643 0208645 0208645 0208645 020864 021084 021084 02257143 02257144 02257144 02257144 02257144 02257144 02257144 02257144 02257144 022577144 02257144 02257144 0225714444 02257144444 02257144444444444444444444444444444444444	were	045 0.3473.99 0.3473.99 0.2473.98 0.27886 0.27881 0.028631 0.173881 0.173881 0.173881 0.173881 0.173881 0.17381 0.17381 0.173821 0.113923 0.113933 0.1139320	ncoun	04b 0.11428.57 0.11428.57 0.1252978 0.1252978 0.1252978 0.1252978 0.1252978 0.1252978 0.12756190
04a 0088834 00088834 0000441 0000441 00185459 00185441 0018653 0018653 0018653 0019563 00115563 0019563 0019563 0019563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00115563 00000000000000000000000000000000000	g this,	g this, ^{du} th	ever e	044 0.2228374 0.2228374 0.2238370 0.223836 0.0103641 0.430147 0.0103588 0.0103588 0.0103588 0.0103588 0.0103588 0.0103588 0.011765 0.0213459 0.0221459 0.0221459000000000000000000000000
026 0.224265 0.224265 0.124265 0.124265 0.124265 0.144264 0.12425 0.1241071 0.014741 0.014722 0.214071 0.014772 0.0147724 0.0106618 0.0106618	readin	0.26 0.26(1029 0.26(1029 0.08(1239) 0.08(1239) 0.08(4559 0.08(4559) 0.08(459)0000000000000000000000000	/e you	03c 0.301471 0.301471 0.301471 0.039063 0.039063 0.039063 0.039063 0.039063 0.0375619 0.0375619 0.0375619 0.0375614 0.0375214000000000000000000000000000000000000
0.05% 0.025735 0.025735 0.025735 0.025735 0.0256731 0.0256731 0.0256731 0.026735 0.018675 0.018675 0.018675 0.027575 0.0257575 0.02757575 0.02757575 0.02757575 0.02757575 0.02757575 0.02757575 0.02757575 0.0275755 0.02757575 0.0275755 0.0275555 0.0275555 0.02755555 0.02755555 0.02755555 0.02755555555555555555555555555555555555	efore	0.3b 0.011029 0.011039 0.01655 0.0655 0.055831 0.053831 0.054959 0.010995 0.010995 0.010995 0.010995 0.010995 0.010995 0.010995 0.0109618 0.006018	2: Hav	036 0.555147 0.555147 0.555147 0.7561713 0.776178 0.776175 0.772056 0.772056 0.772056 0.772057 0.672779 0.556067 0.556067 0.556255 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.557550 0.5575500 0.55750000000000
0.34 0.189955 0.0188955 0.0286731 0.0266731 0.044225 0.044225 0.014525 0.014525 0.014525 0.004525 0.004525	D.1: B	0.34 0.59276 0.59276 0.51723 0.517231 0.517731 0.531731 0.531731 0.531731 0.531731 0.531731 0.531732 0.531732 0.531732 0.531733 0.5	ble D.	034 0.113122 0.113122 0.113123 0.138429 0.138429 0.138429 0.188429 0.188429 0.188429 0.188429 0.188429 0.238577 0.238577 0.238577 0.238577 0.238577 0.238577 0.238577
002 009265 0093255 009255 0081731 01255 0000000000	Table	02c 0.073529 0.073529 0.073529 0.073529 0.073529 0.0752581 0.0753581 0.0753589 0.0753589 0.0753589 0.0753589 0.077135589 0.077135580 0.077135580 0.077135580 0.077135580 0.077135580 0.077135580 0.077135580 0.077155580 0.077155580 0.077155580 0.07755580 0.07755580 0.07755580 0.0775580 0.0775580 0.0775580 0.0775580 0.07755800000000000000000000000000000000	Та	02c 0.257353 0.257353 0.257358 0.257358 0.256739 0.1257368 0.0039965 0.0039965 0.177206 0.177006 0.177206 0.177
02b 092555 0092555 00092555 00092555 00185095 00185055 001850 001850 0000000000		0.176471 0.176471 0.176471 0.176471 0.102341 0.165441 0.165441 0.165441 0.1023841 0.1023841 0.1023841 0.1023841 0.1023842 0.1023842 0.1023842 0.1334592 0.134545 0.134556 0.134556 0.134556 0.134556 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.1355566 0.135556666666666666666666666666666666666		026 0.031142 0.031142 0.0283088 0.1283871 0.1183871 0.1183673575757575757575757575757575757575757
01b 0009265 0009265 0008055 0008029585 000807 000807 000807 0001765 000175 0000000000		01b 0.176471 0.075329 0.075329 0.075329 0.0756471 0.0756471 0.076471 0.07647339 0.0261623 0.0261623 0.0261623 0.0261623 0.0261623 0.0261623 0.0212418 0.0212418 0.0212418 0.021148 0.001148 0.00		01b 0 000 0 0001142 0.0031142 0.0031733 0.0031732 0.0031722 0.014020 0.014020 0.006302 0.0000000000000000000000000000000000
0.1b) (Flavitubk Psprodosis) 0.2b) (Direction of charger in Direction) 0.2b) (Validation of Capetiments) 2b) (Validation of Capetiments)		Oh (Tak inhe hypothesis) (2) (Tak inhe hypothesis) (2) (Christian of charact more) (2) (Careatt more) (2) (C		Oli (Fish finht byochesis) (2) (Fish finht byochesis) (2) (Christian of charge in byochesis) (3) (Christian of charge in byochesis) (3) (Christian in the state in the state (3) (Christian in the state) (3) (Christian of the state) (3) (Christian of the state) (4) (Chr

Table D.3: How severe do you think this problem is?

20b 0.176471 0.117647 0.0126471 0.25624 0.25624706 0.25252525252525252525252525252525252525				
20a 0.117647 0.018824 0.014118 0.044118 0.044118 0.044118 0.05422 0.1143382 0.1143382 0.1143382 0.1143382 0.1143382 0.1143382 0.0153822 0.005022 0.005082 0.00508200000000000000000000000000000000			5 11/16 11/16 11/16 11/17 11/1	
154 0.036765 0.0256765 0.025653 0.1255 0.1255 0.1255 0.012859 0.01285 0.012859 0.012850 0.0128500 0.0128500 0.0128500000000000000000000000000000000000			a a 200 99792 0.14 99792 0.14 99792 0.14 99792 0.14 99792 0.14 14176 0.14 141	
14b 10.088235 0.14769 0.038462 0.038462 0.038462 0.0025 0.0025 0.0025 0.441176 0.441176 0.441176 0.0025 0			228350 200000 228350 228350 200000 228350 228350 200000 228350 228350 200000 228350 200000 228350 200000 200000 200000 200000 200000 200000 2000000	
13a 0.0022576 0.0022548 0.0022548 0.0022548 0.0135489 0.0135461 0.0135461 0.0135461 0.0135461 0.0135459 0.0135459 0.0135595 0.0153595 0.0153595		44 44 1115546 1115545 1115545 1115545 1115655 1115655 1115655 1115655 1115655 1115655 1115655 1113605 1110505 1110505 1110505 1110505 1110505		
12b (0.255294 (0.255294 (0.257159) (0.1377059) (0.1477059) (0.1477059) (0.1477059) (0.1270559) (0.17059) (0.228817) (0.228825) (0.22885) (134 134 0113405 10113405 10113405 1011345 1011345 10117345 10117345 10117345 0010753 0010753 0010012 0010012 0010012 0010012 0010012 0010012 0010025 0010025 0010020 0010020 0010020 0010020 0010020 0010020 0010000 000000 000000 000000 000000 000000	nts?
11b 0.001755 0.047765 0.047705 0.047705 0.04755 0.0475 0.04755 0.04755 0.04755 0.04755 0.04755 0.04755 0.04755 0.04750 0.0025 0.04750 0.04750 0.047706	5	12b 0.00000 0.01720 0.01720 0.01720 0.01720 0.14795 0.14795 0.14795 0.14795 0.14795 0.14795 0.14705 0.	erime	
00a 00a4202 00034622 0004202 0107143 0107143 01014386 01042857 0119048 0119048 0119048 0119048 0119048 0119048 011948 011948 011948 011948 011948 011948 011948 011948 011948 0119587 0117857 0117957 0117857 0117857 0117957 0117857 00000000000000000000000000000000000	oblem		111b 0.02553 0.012756 0.01276 0.01276 0.1776 0.1776 0.17763 0.17763 0.17763 0.05500 0.05500 0.05500 0.012553 0.012553 0.012000 0.012553 0.012000 0.012553 0.012755 0.012553 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.012755 0.0125555 0.0125555 0.0125555 0.0125555 0.0125555 0.01255555555555555555555555555555555555	ng exp
064 11.173669 20.11956 20.11956 20.251905 20.2	his pr		094 00000000000000000000000000000000000	runniı
84 117647 117677 117677 117677 1176777 1176777 1176777 1176777 117677777777	olves t	064 0.257710 0.257710 0.17582 0.17582 0.17582 0.17582 0.17582 0.17582 0.17582 0.17582 0.125877 0.18602 0.142867 0.144867 0.142867 0.142867 0.142867 0.142867 0.142867 0.142867 0.142867	when	
4, 0.025735 0.025735 0.025735 0.025735 0.025735 0.025735 0.025735 0.025735 0.025735 0.025735 0.02575 0.03755 0.03755 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0525 0.048611 0.0575 0.048611 0.0575 0.048611 0.0575 0.048610 0.0575 0.0570 0.0505 0.0500 0.0505 0.0500 0.0505 0.0500 0.0500 0.0505 0.05000 0.05000 0.05000 0.0500000000	tion so		054 0.054 0.02536 0.02536 0.025359 0.025359 0.025359 0.025359 0.025360 0.025360 0.025360 0.025360 0.025360 0.025360 0.025360 0.025360 0.025360 0.02560 0.02600 0.02600 0.02600 0.02600 0.02600 0.02600 0.02600 0.02600 0.02600 0.026000 0.0260000000000	er you
45 11765 1325 1325 1325 1325 1325 1325 1325 132	is solu		046 046 0.003384 0.0033844 0.225964 0.225964 0.225964 0.225964 0.230935 0.003125 0.003125 0.003125 0.003125 0.003125 0.003132 0.003533 0.0031325 0.0031325 0.0031325 0.0031325 0.0031325 0.0031325 0.0031325 0.0031325 0.0031325 0.0035333 0.0035333 0.0035333 0.0035333 0.0035330 0.0035330 0.0035	· hinde
a 3232341 2341176 402178 402178 402775 40275 40075 40000000000	ink th		04b 0.5462 0.54625 0.54625 0.54107 0.54107 0.54107 0.64103 0.64103 0.64103 0.64103 0.64103 0.64103 0.64103 0.64103 0.66803 0.56805 0.56805 0.56805 0.56805 0.56805 0.56805 0.56805 0.57731 0.27731 0.27731 0.17647	ielp or
625 625 625 625 625 625 625 625	you th	: Do you th	044 0.55017000000000000000000000000000000000	tion h
8235 8735 8735 8735 8735 8735 8745 8745 9417 9417 9417 9417 9417 9417 9417 9417	: Do y		0.3c 0.194812 0.194812 0.19482 0.19482 0.19482 0.19463 0.19463 0.21094 0.21094 0.21094 0.21094 0.21094 0.211889 0.2118848 0.2118848 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.211884 0.221855555555555555555555555555555555555	is solu
4 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	e D.4		03b 03b 040935 040935 040935 040935 040954 04074400000000	'ill th
0.34 0.0457 0.10457 0.10038 0.10038 0.01041 0.01457 0.0128577 0.0128577 0.0128577 0.0128577 0.01285777 0.01285777777777777777777777777777777777777	Tabl		034 019457 019457 019458 0200005 018759 018759 041785 022296 012822 024307 01282 012822 021775 02175 0200000000000000000000000000000000000	.5: W
02c 0.101765 0.1012941 0.2211558 0.2211558 0.2211558 0.2211555 0.22558 0.1255 0.13743 0.15 0.138889 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.1397039 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138889 0.138780 0.1387800 0.1387800 0.1387800 0.1387800 0.1387800000000000000000000000000000000000			026 0.33662 0.33662 0.336562 0.00000 0.00000 0.00000 0.23875 0.238756 0.238	ole D
02b 0.058824 0.058824 0.058874 0.102841 0.102841 0.103894 0.147059 0.014705900000000000000000000000000000000000			02b 000000 0100000 0100000 0100000 0100000 0100000 0100000 0100000 0100000 000000	Tał
01b 0.058824 0.058824 0.0497765 0.0497765 0.0497765 0.0487755 0.0487755 0.038824 0.0137665 0.0127675 0.012775 0.0127750 0.012775 0.012775 0.012775 0.012775 0.012750 0.012775 0.0127500000000000000000000000000000000000			016 016 0.00000 0.00000 0.00000 0.00000 0.00000 0.00110 0.000000	
 Hakitabk bypothesis) Ditt Orbitabk bypothesis) Ditt Ontoreal advances of change in hypothesis) Dist Changelian metrics and advances and advances and advances and advances and advances are advances and advances are advances and advances are advances are advances and advances are advances are			11) F Fash indek "Aproduciss) (1) F Fash indek "Aproduciss) (1) F Fash indek "Aproduciss) (1) Constraint metricuss (1) F Canadian metricuss) (1) F Canadian (1) Metricuss) (1) F Canadi	

D.2 Mann–Whitney U test

For this work, a Mann–Whitney U test was performed. See the shorter version in section 6.2. Each bold value is a significant result (< 0.05) with a large effect size.

09.a (Nurview) 25974620 25974620 259974620 25987620 2598762 215516 2110688 2110688 21106882 21106882 2110682 210082 210082 210082 210082 210082 210082 210082 210082 210082 210082 210		09.a (Nur 0.02609) 0.0240515 0.0240515 0.0240515 0.034584 0.034584 0.034585 0.034585 0.034585 0.034585 0.0396805 0.0396805 0.039853 0.03580500000000000000000000000000000000		094 (Nur 0.7256)-10 0.574856 0.574856 0.574856 0.574856 0.034440 0.276956 0.277956 0.2769566 0.2769566 0.2769566 0.2769566 0.2769566 0.2769566 0.276956 0.276956 0.276956 0.27	
664.(Failure checks) 0.054.(Failure checks) 0.0310250 0.0382050 0.038205 0.038205 0.012057 0.012057 0.012051 0.012051 0.020430 0.020400000000000000000000000000000000		066. (Failure checks) 0.64525 0.045031 0.045031 0.045031 0.045031 0.045031 0.045031 0.045031 0.045035 0.05515 0.005112 0.005136 0.005256 0.005256 0.005256 0.005256 0.0057556 0.0056575 0.0056575 0.0056575 0.005756 0.005756 0.005756 0.005756 0.005675 0.005756756 0.005756 0.005756 0.005756 0.005756 0.005756 0.005756 0.0		066 (Failure checks) 1 016 (99 16) 99 1 1 272 1 1 272 1 1 2018 1 1 2018 1 1 2018 1 1 2018 1 1 2018 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2	
ous experiments)		ous cyperiments)		De 2. Mann–Whitne	y U test
054 (Simultan 0054 (Simultan 0.0058447 0.0058447 0.0186414 0.018644 0.014843 0.014843 0.014843 0.014843 0.014830 0.013853 0.01385447 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.0138547 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.00385567 0.000385567 0.00385567 0.00385567 0.00385567		05.4 (Simulan 0.14,3998 0.14,3998 0.14,3998 0.14,3998 0.14,3998 0.17,3738 0.17,3728 0,		05.4 (Shruhtan 0.33721/201833 0.337212 0.248701 0.248701 0.29512 0.290512 0.290512 0.29077 0.248702 0.248702 0.2488702 0.2488702 0.2488702 0.2488702 0.2488702 0.248840 0.2488702 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.2488400 0.24884000000000000000000000000000000000	
04. (Wribadding reaults) 04. (Wribadding reaults) 0.0011153 0.0380123 0.0311546 0.0131546 0.0131546 0.0131546 0.0131548 0.013548 0.0136548 0.0137548 0.013655548 0.0136555555555555555555555555555555555555		04. (Withbulling results) 0.539477 002 0.539477 002 0.531437 0.0027345 0.0027345 0.0027345 0.0027345 0.0027345 0.331882 0.331882 0.331882 0.331882 0.331882 0.331882 0.331882 0.03258 0.0155855 0.0155855 0.0155855 0.015585555 0.015555555555555555555555555555555		04. (Winbidding results) 2237555 0.056531 0.056531 0.057148 0.0071482 0.03140000000000000000000000000000000000	
0.0) (Calculate minimum duration) 0.005 (Calculate minimum duration) 0.005 (SS) 0.005 (SS) 0.005 (SS) 0.005 (SS) 0.012 (SS) 0.014 (S		04b (Calaba later minimum duration) 0.030753 0.030753 0.030653 0.0164536 0.0164563 0.0164863 0.0168463 0.0168463 0.036463 0.0364645 0.036465 0.036656 0.036656 0.036656 0.036656 0.036656 0.036656 0.03665656 0.03665656 0.03665656 0.03665656 0.03665656 0.03665656565656 0.0366565656565656565656565656565656565656		04b (Calculate minimum duration) 0448.570 0448.570 0.0506.173 0.0506.173 0.0507.174 0.0507.174 0.0507.174 0.0507.174 0.0507.174 0.0717.855 0.0717.855 0.0717.855 0.0717.855 0.0211.12 0.0717.855 0.0211.12 0.0717.855 0.0211.12 0.	
044, (Minimum effect 5126) 0144, (Minimum effect 5126) 0124545 0124545 0124545 0124545 0124545 0124545 0140047 0140048 0140048 0116426 0106456 010666 010666 0106666 0106666 0106666 0106666 0106666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 01066666 00000000	ı existed?	044, (Minimum effect 3420) 002834/57 002834/57 002834/57 002834/57 002834/57 003834/57 003834/57 003834/57 003834/57 003834/57 003834/57 003424/1 0000124/1 00000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 0000124/1 000000000000000000000000000000000000	rself?	044, 054, 054, 054, 054, 054, 054, 054,	2
0.154 (Churning users) 0.154 (Churning users) 0.154 (Churning users) 0.154 (Churning users) 0.154 (Churning users) 0.157 (Churning users) 0.156 (Churning users) 0.156 (Churning users) 0.171 (Churning users)	iis problem	0.5 c (C)mming users) 0.05 mming users) 0.059078 0.059078 0.059078 0.059078 0.059078 0.059078 0.05808 0.05808 0.05808 0.05808 0.05808 0.05808 0.05808 0.05808 0.05756 0.047567 0.067756 0.067556 0.067556 0.067556 0.067556 0.067556 0.067556 0.067556 0.067556 0.067556 0.067556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.065556 0.075556 0.075556 0.075556 0.075556 0.075556 0.075556 0.075556 0.075556 0.075556 0.075556 0.077556 0.075556 0.075556 0.075556 0.07555656 0.0755566 0.0755566 0.07	oblem you	0.15 c (Churning users) 0.134764 0.10025 0.10025 0.10025 0.2803964 0.2803964 0.280397 0.589347 0.589347 0.589347 0.589347 0.237056 0.237056 0.237105 0.237105 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239484 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.239588 0.2395888 0.239588 0.239588 0.23958	problem is
0.08 (Competitor atley) 0.089 (Competitor atley) 0.080 (Competitor atley) 0.048623 0.048623 0.028534 0.029344 0.041742 0.041742 0.041742 0.041742 0.041742 0.041742 0.0139346 0.0213542 0.021354 0.00144555 0.00144555555555555555555555555555555555	ou aware tl	00b (Competitor sulty) 0.957017 0.957017 0.957017 0.977775 0.977775 0.977775 0.977775 0.977775 0.977775 0.97775 0.97775 0.07438 0.07438 0.07438 0.07438 0.07438 0.003577 0.003577 0.003577 0.003577	red this pr	CB) (C) mpetitor safety) (C) (C) mpetitor safety) (C) (C) (C) (C) (C) (C) (C) (C) (C) (C)	think this
(Ba. (Technical debt) (Ba. (Technical debt) (1931) (1932)	iis, were y	03, (Tischnisal debt) 001255 001255 00127697 00127697 00127697 00127695 0012755 0012755 00127555 00127555 0011100 0011100 00127155 00127155 00127155 00127155	er encounte	03. (The Antical debt) 0.399644 0.399644 0.399644 0.0325612 0.039644 0.0394774 0.394774 0.394774 0.394774 0.394773 0.394774 0.394773 0.3947773 0.3947777 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.3947773 0.39477473 0.3947747474747474747474747474747474747474	ere do you
Cl2; (Caurdrail metrics) 0.25; (Caurdrail metrics) 0.65; 12038 0.65(0)33 0.65(0)33 0.65(0)33 0.48625 0.186641 0.186641 0.186641 0.186641 0.1866400 0.1866400 0.1866400 0.1866400 0.1866400 0.1866400 0.1866400 0.1866400000000000000000000000000000000000	e reading th	025 (Gaurdraft metrics) 0254137 0254137 02551381 02551432 02551432 02551432 02551432 02551432 02551432 0255153 0255151 0265511 0265511 0265511 0265511 0265511 0265513 026555 0217556 0277556 0277556 0277556 0277556 0277556 0277556	ave you eve	02. (Gaundrall metrics) 0. (Standrall metrics	8: How sev
0.00 Direction of change in hypothesis) 0.00 Direction of change in hypothesis) 0.0978388 0.072688 0.072688 0.072688 0.0717478 0.0073179 0.0073173 0.0073173 0.0053173	able D.6: Befor	CDb (Direction of change in hypothesis) 0.0524003 0.0524003 0.052403 0.052403 0.052403 0.052403 0.053403 0.054403 0.054403 0.054403 0.054403 0.054403 0.054403 0.014504 0.014514 0.015414 0.012414 0.012414	Table D.7: H	CDA (Direction of change in hypothesis) 0.88744 0.088744 0.086874 0.086845 0.002835 0.002835 0.002835 0.0028345 0.0028345 0.0028345 0.002845 0.002845 0.057465 0.057465 0.05505 0.05505	Table D.
018 (Falaintako hypothesis) 018 (Falaintako hypothesis) 018 (Falaintako hypothesis) 018 (Falaintako 108 (Falaintako 108 (Falainta) 018 (Falai	L	01b (Fa)JS1 (IIIAble by pedhesis) 0.9592031 0.260031 0.260031 0.260031 0.260031 0.213345 0.213345 0.213345 0.213345 0.213345 0.213345 0.213345 0.223345 0.23345 0.233455 0.233455 0.233455 0.233455 0.2334555 0.2334555 0.23345555 0.2334555556 0.2334555556 0.233455556 0.233455556556 0.23345555655656565656565656565656565656565		01b (Falsifitible by pothesis) 0.0858744 0.0838744 0.0242548 0.0242548 0.0145645 0.0145644 0.0146549 0.0146592 0.0167992 0.0167992 0.016792 0.016792 0.017525 0.0175648 0.02555 0.027553 0.027550 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.027553 0.0275553 0.0275553 0.0275553 0.02755555555555555555555555555555555555	
01h Flakitabe hypothesia) 02h Draftabe hypothesia) 02h Draftabe draftage in by pothesia) 02h Draftabel data) 02h Channel data atzl 02h Channel datzl 02h Channel data atzl 02h Channel data atzl 02h Cha		01b (Flakithe hypothesis) 02b (Flakithe hypothesis) 02b (Discretion of change in hypothesis) 02b (Discretion at device) 02b (Competitive atley) 02b (Competitive atley) 04b (Change muture) 04b (Mattumm effect ster) 04b (Mattum effect ster) 04b (Mattum effect ster) 04b (Mattum effect ster) 15b (Mattum effect st		01b (Flakithek bypothesis) 02b (Charection of change in hypothesis) 02b (Charection of change in hypothesis) 03b (Competition and the second and the seco	

094 (Num 0.777820 0.777820 0.777820 0.77952 0.159525 0.159525 0.0125887 0.0125887 0.0125887 0.0125887 0.0119141 0.0751051 0.075100 0.075100 0.075100 0.075100 0.075100 0.07510000000000000000000000000000000000		0a (Num 1 0351356 0351356 035878 0.35878 0.35878 0.358738 0.3587348 0.3587348 0.3587348 0.3587348 0.3587348 0.357348 0.256348 0.257348 0.257458 0.2	
064. (Faline check) 064. (Faline check) 0.050231 0.050231 0.050031 0.050035 0.045031 0.045031 0.043737 0.0438235 0.0438235 0.0438235 0.0446647 0.044647 0.04467 0.044647 0.04467 0.044677 0.044647 0.044647 0.0446477 0.0446477 0.04464777000000000000000000000000000000		046. (Failure checks) 046. (Failure checks) 053753 05175345 01374455 01374455 01374455 01374455 01374455 0137455 0137550 0137550 01000000000000000000000000000000000	
D. STATISTIC TABLES		us experiments)	
64, (Shimultuneon 0, 0, (Shimultuneon 0, 0, 0, 15) 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0		054 (Simultaneo 053 (Simultaneo 0.7907054 0.87721 0.807721 0.15721 0.15721 0.15721 0.15721 0.15721 0.256054 0.256054 0.256054 0.256054 0.256054 0.256054 0.256054 0.256057 0.257057 0.256057 0.257057 0.256057 0.257057 0.256057 0.257057 0.2	
94, (Withoding realis) 94, (Withoding realis) 0.059799 0.059799 0.0393947 0.039347 0.031346 0.031346 0.031346 0.031345 0.0323432 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.032342 0.03242 0.032420 0.032420 0.03242000000000000000000000000000000000		04c (Whibding results) 0.049266 0.049266 0.049266 0.049215 0.049215 0.049215 0.0492215 0.0492215 0.059214 0.0592361 0.0592264 0.0252763 0.0252267 0.0252267 0.0252267	
(4) (Clasher minimum duration) (0014-0) (0014-0) (0014-0) (0014-0) (00104-0)		04b (Calculate minimum daration) 0000345 0000345 0000345 0000375 0000375 0000375 0000355 00000355 000035 0000005 00000000	
44. (Mainimm effect size) (0.038952 0.038952 0.0119926 0.0119926 0.0119926 0.0119926 0.011926 0.011926 0.011926 0.012588 0.014323 0.0143333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.014333 0.0143333 0.0143333 0.014333333 0.014333333 0.01433333 0.01433333 0.01	em?	044 (Minimum effect size) 0006613 01006613 0101048 0101044 01000445 01000454 01000454 0100054 0100054 0100054 0100054 0100054 0100054 0100054 0100054 0100054 0100054 0100054 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 0100055 010055 00055005 000055 00055005 0005005	experiments?
 0. 0.4.5.111 0.4.6.5.111 0.4.6.5.111 0.4.6.5.111 0.0.822.54 0.0.922.54 	s this probl	0 Gbc (Churning users) 0.443-964 0.443-964 0.443-964 0.9443-257 0.9443-258 0.9443-258 0.9443-258 0.9443-258 0.0025-258 0.0258-258 0.9559-258 0.9559-258 0.9559-258 0.9559-258 0.0251-34 0.0252-34 0.	en running
0.001153 0.001153 0.011153 0.011153 0.011153 0.011153 0.011153 0.01153 0.001153 0.001153 0.00001153 0.00001153 0.000010000000000000000000000000000000	ak this solution solves	03b (Conjuctine sufet 0.051363 0.053136 0.056378078000000000000000000000000	ler you wh
 (b) (b) (b) (b) (b) (b) (b) (b) (b) (b)		0 (34 (Trebnical deb) 0.85(4) 0.95(4) 0.95(4) 0.976(9) 0.976(9) 0.976(9) 0.97136 0.97136 0.00014 0.00014 0.025(18) 0.025	aelp or hind
 COA, (Carandral) Interfic- (24, 40007 COA, (Carandral) Interfic- (24, 40007 CAA, 0007 <	Do you thi	CCc (Characteril Ineutic CCc (Characteril Ineutic CCC (Characteril CCC (Characteril CCC (Characteril CCC (Characteril CCC (CCC) CCC (CCC) CCCC) CCC (CCC) CCC (CCC) CCC (CCC) CCC (CCC) CCC (CCC) CCC (CCC) CCC (CCC) CCC)	s solution h
0.0 fblaction of change in by polatesis 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.73 4307 0.003 6307 0.003 6307 0.00307 0.003 6307 0.00307 0.00307 0.00307 0.0030000000000	Table D.9:	CDb (Direction of change in hypothesis 0.95 (2001) 0.95 (2001) 0.05 (2001) 0.05 (2001) 0.03 (2001) 0.0	le D.10: Will thi
0(h, fralutitable by porthesis) 0(h, fralutitable by porthesis) 0.374.308 0.374.308 0.343427 0.343427 0.013929 0.013929 0.013929 0.013929 0.013929 0.0139213 0.0139213 0.0397140		01b (Falsifiable by podtesis) 1038/744 1038/744 1038/748 04318/64 04318/64 04318/64 04318/65 04318/65 04308/65 04308/65 04308/65 04408/14 04408	Tab
1) () () () () () () () () () () () () ()		01) h (Fakifabb bypothesis) 02) Contexcination doubling in hypothesis) 02) Contexcination doubling in hypothesis) 03) (Competitive address) 03) (Competitive address) 03) (Competitive address) 04) (Minimum direct address) 04) (Minimum direct address) 04) (Minimum direct address) 04) (Minimum direct address) 05) (Minimum direct address) 13) (Minimum direct address) 14) (Minimum direct address) 14) (Minimum direct address) 14) (Minimum direct address) 15) (Minimum direct address) 16) (Minimum direc	