

DELFT UNIVERSITY OF TECHNOLOGY

PRE-REGISTRATION REPORT

Costs of predictions in hate speech detection

Authors:

Philippe Lammerts

Prof. dr. ir. G.J.P.M. Houben, thesis advisor

Dr. J. Yang, TU Delft, daily supervisor

Dr. Y-C. Hsu, TU Delft, co-daily supervisor

P. Lippmann, TU Delft, co-daily supervisor

May 12, 2022



Abstract

This document is prepared for the Human Research Ethics Committee review at TU Delft. It describes the plan for an experiment for the Thesis project called "Building a smart rejector for detecting hate speech". This document follows the pre-registration plan suggested by [20]. This document elaborates on the goal, the exact procedure, and the analysis of the experiment.

In this experiment, we will conduct an anonymous survey where we ask human subjects to judge different scenarios of hate speech detection on social media. The goal of this experiment is to find out the relative cost values between these scenarios. We will use the results in the Thesis project to build a smart rejector for hate speech detection.

Contents

1	Research Question and Hypotheses	2
2	Method	3
2.1	Background	3
2.1.1	Likert	3
2.1.2	Magnitude Estimation	3
2.2	Design	4
2.2.1	Idea	4
2.3	Variables	4
2.3.1	Independent variables	4
2.3.2	Confounding variables	5
2.3.3	Control variables	5
2.3.4	Dependent variables	5
2.4	Planned sample	5
2.4.1	Participant Inclusion and Exclusion Criteria	6
2.4.2	Participant Compensation	6
2.4.3	Sample size	6
2.5	Materials	6
2.5.1	Survey tool	6
2.5.2	Data	6
3	Procedure	7
3.1	Survey	7
3.2	Example FN scenario with ME scale	8
3.3	Example FP scenario with 100-level scale	8
3.4	Example rejection scenario with 100-level scale	9
4	Analysis	9
4.1	Costs	9
4.2	Reliability	9
4.3	Validity	10
A	Presentation texts	10
A.1	Consent	10
A.2	Short introduction ME-100	11
A.3	Short introduction 100-ME	11
A.4	Introduction	11
A.5	100-level scale explanation	11
A.6	ME scale explanation (inspired by [14])	12
A.7	Attention check	12
A.8	Training phase ME	12
A.9	ME calibration	12

1 Research Question and Hypotheses

The amount of hateful content that is spread online on social media platforms remains a serious problem. Manual moderation is still the most reliable solution but is simply infeasible due to the large amount of data generated every second on social media platforms [2]. There exist automated solutions for detecting hate speech, and most of these use Machine Learning models. However, these models tend to be unreliable as they often perform poor on deployment data [2, 8]. One study found that the F1 scores reduce significantly (69% F1 score drop in the worst case) when training a hate speech detection model on one dataset and evaluating it using another dataset [8].

Therefore in this project, we focus on Machine Learning (ML) models with a reject option. The goal of the reject option is to reject a prediction when the model is not confident enough [9]. This Thesis project is about building a smart rejector for detecting hate speech. A system in which the machine assists the human in detecting hate speech automatically and in which the human makes the decisions when the machine is not confident enough.

The first goal of the project was to find a confidence metric that calculates the optimal rejection threshold. This threshold is calculated based on the confidence values of the ML predictions and a set of cost values. These cost values represent the impact of the True Positive (TP: predicting content correctly as hateful), True Negative (TN: correctly predicting content as non-hateful), False Positive (FP: predicting content as hateful while it is non-hateful), False Negative (FN: predicting content as non-hateful while it is hateful) and rejected predictions. Rejecting a prediction means that the ML predictor was not confident enough and, therefore, a human moderator needs to make the final judgment. We understand that there are different ways of handling rejected predictions. In this project, we treat rejected predictions as content that remains publicly visible online for, at most, 24 hours until a human moderator decides to remove/tolerate it. We based the timespan of 24 hours on a news article that stated that Germany fines social media platforms if they do not remove hate speech within 24 hours¹.

The second step of this project is to find out how we can retrieve the relative cost values of TP, TN, FP, FN, and rejected predictions in hate speech detection. By relative costs, we mean to figure out, for example, the cost ratio between an FP and an FN prediction. Expressing these cost values in money spent/saved is infeasible in hate speech detection due to many uncertainties. So in this experiment, we aim to define the cost values in a subjective manner by analyzing the subjects' opinions about different hate speech detection scenarios. The main research question and our hypotheses of this experiment are as follows:

[RQ]: How can we determine the relative costs of rejections and True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) predictions in hate speech detection?

- We hypothesize that the cost of an FN is greater than the cost of an FP. We believe that allowing hateful content to be publicly visible has a more negative impact on social media users than filtering out neutral content. Therefore, we think that the cost of an FN is greater than the cost of an FP.
- We hypothesize that the gain of an TP is greater than the gain of an TN. We believe that correctly predicting hateful content is more valuable to social media users than correctly predicting non-hateful content. Therefore, we think that the gain of an TP is greater than the gain of an TN.
- We hypothesize that the cost of a rejection is lower than the average cost of an FP and an FN prediction. The key assumption of using ML models with a reject option is that the cost of a rejection should always be lower than the cost of an incorrect decision.
- We hypothesize that Magnitude Estimation (ME) is a suitable technique for retrieving the relative costs. For details on ME refer to Section 2. ME seems like a promising technique for retrieving ratio data from judgements about hate speech detection scenarios. We use a 100-level numerical scale for validation. We expect that both scales are correlated and will give similar judgements. Although we also expect the 100-level scale to be suitable for retrieving opinions about the different hate speech detection scenarios, it does not provide the ratio data we need. We also expect that the inter-rater reliability for the 100-level scale will be higher than for the ME scale since the ME scale provides more response freedom. We also expect this since the authors of [18] concluded that the inter-rater reliability of the 100-level scale is higher than the ME scale.

¹<https://www.nytimes.com/2017/06/30/business/germany-facebook-google-twitter.html>

2 Method

In this experiment, we show different scenarios to the human subjects using a survey and ask them whether they agree, disagree or are neutral about the decisions of a fictional social media platform we call SocialNet. These scenarios represent TP, TN, FP, FN, and rejected predictions. For example, we simulate an FN by showing a hateful post to the subject and explain that the detection algorithm of SocialNet did not detect any hate speech. The subject then indicates whether they agree, disagree, or are neutral with this decision. By analyzing the results of a group of subjects, we aim to conclude that an FN scenario is considered, for example, three times worse as an FP scenario. This section first gives some background information and then elaborates on the exact procedure of the experiment.

2.1 Background

Initially, we considered using Likert scales as the response scales of our survey. However, as we will explain in the following subsection, Likert scales are not suitable for retrieving ratio values between the different scenarios of TP, TN, FP, FN, and rejections. Therefore, we will use a technique called Magnitude Estimation in our survey.

2.1.1 Likert

Most papers use Likert scales for retrieving the opinions of a group of subjects. Likert scales are a series of multiple Likert-type questions where subjects can answer questions with several response alternatives [5]. So in our case, we could use a bipolar scale with seven response alternatives ranging from 'strongly disagree' to 'strongly agree', including a 'neutral' midpoint. However, there is a lot of discussion in the literature about how we should analyze these Likert scales [5, 1, 16, 15]. The scale of the questions is ordinal, which means that we do know the ranking of the responses, but we do not have an exact measurement of the distances between the response items [1]. For example, we know that 'strongly agree' is higher in rank than 'agree', but not the exact distance between the two responses and whether it is greater than the distance between the 'neutral' and the 'somewhat agree' responses. Therefore, we cannot use parametric statistics, such as calculating the mean, when analyzing the data [1]. Other papers argue that we can treat a Likert scale that consists of multiple Likert items as interval data and, therefore, applying parametric statistics will not affect the conclusions [5, 16, 15]. So, we can calculate mean scores for TP, TN, FP, FN, and rejection scenarios and compare these with each other. For example, we can then verify that the mean value of FN cases is greater than the mean value of FP cases and conclude that the cost of an FN is greater than the cost of an FP. Analyzing Likert scales from our surveys would at most provide us with interval data (data for which we know the order, and we can measure the distances [1]). However, we need to have ratio data in this project since we want to know the exact ratios between the cost values of the TP, TN, FP, FN, and rejection scenarios.

2.1.2 Magnitude Estimation

In 2.1.1, we concluded the Likert scales are not suitable since they do not provide ratio data. In this research, we want to experiment with a technique called Magnitude Estimation (ME). The ME technique originates from psychophysicists where human subjects need to give a quantitative estimation of sensory magnitudes [19]. For example, in one experiment, human subjects are asked to assign any number that reflects their perception of the loudness of a range of sounds [19]. If the human subjects perceive the succeeding sound as twice as loud, they should assign a number to it that is twice as large. Researchers applied the ME technique to different types of physical stimuli (line length, brightness, duration, etc.) and proved that the results are reproducible and the data has ratio properties [14]. Other works have shown that the ME technique is also useful for more abstract types of stimuli, such as judging the relevance of documents [12], the linguistic acceptability of sentences [3], the relative quantification of political opinions [11], and the usability of system interfaces [13]. Therefore, we think that ME is a promising method for judging the relative costs of scenarios in hate speech detection.

The main advantage of ME is that it provides the ratio scale properties we need. Another advantage is that the scale is unbounded compared to other commonly used response scales, such as Likert scales. For example, suppose we first show a scenario and the subject provides a 'strongly disagree' judgment. Suppose we now present an even worse scenario. The subject is now limited to the response items in the Likert scale and can only give the same 'strongly disagree' judgement. We do not have this problem when using ME because the subject always has the freedom to assign a value of disagreement that is even larger. There are two main drawbacks of using ME. First, we need to normalize the results. Second, we need to validate if we can use ME to measure the subjects' judgements for different scenarios of hate speech detection.

The resulting data needs to be normalized since each subject can use any value they like. For example, one may judge the scenarios using 1, 2, and 10, while another may use values of 100, 200, and 1000. Luckily, there are different solutions for normalizing ME data, such as modulus normalization or external calibration [14]. The most commonly used method for modulus normalization is geometric averaging since this preserves the ratio information [14, 13]. However, as opposed to the unipolar scales used in [3, 13], we are using bipolar scales (disagree-agree). By including 0 (neutral) and negative values (disagree), we cannot use geometric averaging anymore because it uses log calculations [14]. Using the algorithmic mean is also not an option since it would destroy the ratio scale properties [14]. Therefore [14] proposed an external calibration method that keeps the ratio scale properties. This method instructs the subjects at the end of the survey to indicate which verbal label (such as strongly agree) corresponds with which numerical value they used before [14]. Then we calculate the average value, called the pivot value, and divide each ME value by the pivot value [14]. We can now normalize the ME values among all subjects and keep the ratios.

Most papers that use the ME method apply some form of validation. Cross-modality is a technique that is often applied to validate the ME results [3]. Psychophysicists compare the magnitude estimates to the physical stimuli [3]. They analyze the correlation between the magnitude estimates and the physical stimuli by taking the log of each value and plotting them against each other [3]. In the case of estimating line lengths, we can easily vary the line length, for example, by showing a line that is twice as long as the previous line. Subjects can then estimate the line length using a number that is twice as large. However, this becomes more difficult in the social science and psychology domains. In hate speech detection and other applications in social science and psychology, we do not have an exact measure of the stimulus [3]. Luckily, related work has shown that ME is still a suitable technique for eliciting opinions about different types of non-physical stimuli [3, 13, 12, 11]. We can validate the magnitude estimates by adopting the cross-modality technique but instead compare judgements against judgements [3, 11]. We validate our findings by checking the correlation between two different measures of the same thing. We call this a form of convergent validation [6]. Therefore, we need to compare the ME scale against another. When there is a correlation between the two scales, we can conclude that they measure the same thing that the validity of the ME technique has increased.

2.2 Design

We will explain in this chapter the experiment idea and the experimental setup (variables considered and materials used). All the collected data is completely anonymous. We will inform the participants not to put personal identifiers in their answers.

2.2.1 Idea

In this experiment, we will shuffle and present the TP, TN, FP, FN, and rejection settings to the subjects. We will use two scales in our experiments: the ME scale and the 100-level scale. The 100-level scale is a bounded scale that consists of 100 numerical levels. We will use the 100-level scale to validate whether we can use the ME scale to measure the subjects' opinions about the different types of scenarios. This validation is a form of convergent validation where we try two types of measures to see if they measure the same thing [6]. The 100-level scale is easier to understand than ME, does not require normalization, and provides more flexibility than Likert scales [12].

2.3 Variables

We will list all independent, dependent, confounding, and control variables analyzed in our experiment in this section. Important to note is that we do not intend to study the confounding variables demographic and age. The main reason for this is that we do not focus on the effects of subject characteristics on the outcomes but only want to study whether the ME technique is suitable for retrieving opinions about hate speech detection predictions. We left the investigation of the effects of subject characteristics on the outcomes to future research. According to [7], there is no significant difference in how men and women perceive hate. Therefore, we do not consider gender as a confounding variable.

2.3.1 Independent variables

- **Scenarios**

- **True Positive** Show a hateful post to the user and explain that SocialNet detected hate and removed the post.

- **True Negative** Show a non-hateful post to the user and explain that SocialNet did not detect hate and tolerated the post.
- **False Positive** Show a non-hateful post to the user and explain that SocialNet detected hate and removed the post.
- **False Negative** Show a hateful post to the users and explain that SocialNet did not detect hate and tolerated the post.
- **Rejection**
 - * Show a hateful post to the user and explain that SocialNet was uncertain whether the post was hateful or not. An internal moderator will need to check the post within 24 hours. Meanwhile the post remains visible.
 - * Or, show a non-hateful post to the user and explain that SocialNet was uncertain whether the post was hateful or not. An internal moderator will need to check the post within 24 hours. Meanwhile the post remains visible.
- **Social media posts**
 - Neutral content.
 - Hateful content that is either sexist or racist.

2.3.2 Confounding variables

- **Demographic** People from different countries might have different perceptions and definitions of hate speech and how we should deal with it (allowing hateful content or removing it).
- **Age** People of different ages might have different perceptions and definitions of hate speech and how we should deal with it (allowing hateful content or removing it).
- **Learning effect of the scales** Since each subject needs to use both scales, we might introduce a learning effect problem. Once the subject learns how to rate with one scale and then uses another, then the results of the second scale might be affected.

2.3.3 Control variables

- **Scales** Because of our confounding variable learning effect of the scales, we need to cancel the learning effect out. We do this by splitting the group into two and assigning a different order of scales to each group. For example, the first group needs to answer the questions using the ME scale first and then 100-level scale. The second group needs to answer the questions using the 100-level scale first and then the ME scale.
- **Presentation of scenarios** In both experiments, we will present the individual or pairs of scenarios in a random order to reduce bias.
- **Content of the posts** All social media posts are sampled from existing datasets containing both hateful and non-hateful tweets.

2.3.4 Dependent variables

- **Reliability** Measured using the Krippendorff's alpha where values larger than 0.8 indicate reliable conclusions and values larger than 0.6 indicate tentative conclusions [10].
- **Validity** Convergent validity (if two different measures measure the same thing) [6]. Measured by calculating the correlation between the magnitude estimates and the response values from the 100 level scale.
- **Costs of TP, TN, FP, FN, and rejection scenarios** Measured by calculating the mean of the normalized magnitude estimates of each scenario.

2.4 Planned sample

We will explain in this section how we will recruit the participants, give our sample size, and set our stopping and exclusion rules.

2.4.1 Participant Inclusion and Exclusion Criteria

We will use the [Prolific](#) platform for recruiting online participants. We will use the following inclusion criteria for our participants:

- 18 years of age and older since we are showing offensive language in the experiment.
- Fluent in English.
- Approval rating over 90% on the Prolific platform.
- Use one of the following social media platforms regularly (at least once a month): Facebook, Twitter, YouTube, LinkedIn, Pinterest, Google Plus, Tumblr, Instagram, Reddit, VK, Flickr, Vine.co, Meetup, ask.fm, Snapchat, TikTok, Medium.

We use the following exclusion/rejection criteria:

- Participants who fail the attention check. Around the start of the survey, we will include an Instructional Manipulation Check to check if the user pays attention to the survey².
- Participants who do not complete all questions.
- Participants who do not agree with the informed consent before the start of the survey. We are not allowed to collect and process their data if they do not consent.

2.4.2 Participant Compensation

Every participant will be paid based on the hourly wage of 9.0 GBP (about 10,67 Euro), indicated as good pay by the platform³.

2.4.3 Sample size

There are 4.55 billion active social media users⁴. We choose a 95% Confidence Interval (CI) and 10% Margin of Error (MoE) for this study. So for 95% of the time, our observations will fall within a 10% interval [17]. According to [17], we need a sample size of 96 participants to reach the desired CI and MoE values. We chose 10% MoE since we have a limiting budget. We will first conduct a pilot survey for 5 participants to gather feedback and check if we need to improve things before the actual experiment. We want to determine the average workload using the pilot survey and decide whether its possible to reduce the MoE by increasing the number of participants. For the pilot survey, we will use 5 participants. Therefore, in total we will need $5 + 96 = 101$ participants.

2.5 Materials

2.5.1 Survey tool

We use [LimeSurvey](#) as our survey tool since it supports all the features we need, and its (discounted) subscription price is 17 euros per month.

2.5.2 Data

During experiment Type 2, we need to cover all 10 different pairs of scenarios which require 20 social media posts. Furthermore, we need each subject to rate all 10 pairs with each scale. Therefore, we need 40 different social media posts where 20 are hateful and 20 non-hateful. For experiment Type 1, we will use the same posts to create 40 individual scenarios. We used the datasets from [21] and [4] and randomly sampled both neutral and hateful samples. Both datasets contain tweet messages from Twitter that are neutral, racist or sexist. We filtered out the Twitter replies and mentions since the context is not always clear in these messages.

²<https://researcher-help.prolific.co/hc/en-gb/articles/360009223553>

³<https://prolific.co/pricing>

⁴<https://datareportal.com/reports/digital-2021-october-global-statshot>

3 Procedure

This section will explain all steps of the survey. [Appendix A](#) contains the presentation texts. All subjects have 10 seconds before they can continue to the next question. There is no limit to the amount of time they spend on each question.

3.1 Survey

WARNING: the examples used in this section contain content that may make some people feel uncomfortable.

Step 1: provide informed consent

- Show the informed consent (with checkboxes for giving consent).
- Proceed to the next step only for the participants who give consent.

Step 2: introduction

- Show introductory text about what is expected from the subject.
- We split all participants up into two groups.
- The first group first uses the ME scale to rate the first half of all scenarios and then uses the 100-level scale for the second half.
- The second group first uses the 100-level scale to rate the first half of all scenarios and then uses the ME scale for the second half.
- Provide an explanation in the introduction about the first scale.

Step 3: Attention check

- Simple attention check where we ask the subject to select one option.

Step 4a: ME practice phase (when ME is used first)

- To let subjects learn how to use ME, we first run a practice phase where we shuffle and present 5 different line lengths.
- Each subject needs to estimate the line length using any positive value.

Step 4b: first half of all scenarios using the first scale (ME or 100-level)

- Type 1: show 20 different scenarios in random order: 4 TP, 4 TN, 4 FP, 4 FN, and 4 rejection.
- Type 2: show 10 different pairs of scenarios in random order: TP-TN, FP-FN, TN-FN, TP-FP, TP-FN, TN-FP, FP-rejection, FN-rejection, TP-rejection, and TN-rejection.

Step 4c: ME calibration phase (when ME is used first)

- Ask subjects what numbers correspond with the following verbal labels: Strongly disagree, disagree, somewhat disagree, somewhat agree, agree, strongly agree.
- We will use this to normalize the magnitude estimates as explained in the external calibration method from [\[14\]](#).

Step 5a: ME practice phase (when ME is used second)

- Explain that we are going to use the ME scale in the second part.
- To let subjects learn how to use ME, we first run a practice phase where we shuffle and present 5 different line lengths.
- Each subject needs to estimate the line length using any positive value.

Step 5b: second half of all scenarios using the second scale (ME or 100-level)

- In case we now use the 100-level scale, explain that we now switch to the 100-level scale and how this scale works.
- Type 1: show 20 different scenarios in random order: 4 TP, 4 TN, 4 FP, 4 FN, and 4 rejection

- Type 2: show 10 different pairs of scenarios in random order: TP-TN, FP-FN, TN-FN, TP-FP, TP-FN, TN-FP, FP-rejection, FN-rejection, TP-rejection, and TN-rejection

Step 5c: ME calibration phase (when ME is used second)

- Ask subjects what numbers correspond with the following verbal labels: Strongly disagree, disagree, somewhat disagree, somewhat agree, agree, strongly agree.
- We will use this to normalize the magnitude estimates as explained in the external calibration method from [14].

Step 6: finish

- Show a thank you message and redirect the users to Prolific to complete the task.

3.2 Example FN scenario with ME scale

Suppose we have the following post:

I really hate all immigrants, they steal our homes!!!

Question 1: Please indicate **how hateful** you feel this post is using any number greater than or equal to 0. A large number means you think this post is very hateful, and a small number means you think its slightly hateful. If you feel that this post is not hateful at all, enter 0.

****Show ME input field that allows all positive values greater than 0.****

Question 2: SocialNet **removed** the post. The detection system labeled this post as hateful. And so, SocialNet removed the post immediately.

Please indicate whether you agree, disagree, or are neutral about SocialNets decision.

Question 3: Please indicate how much you agree or disagree with SocialNets decision using any positive number. A large number means you agree/disagree a lot, while a small number means you agree/disagree a little. If you feel neutral about SocialNets decision, select neutral in the field above.

****Show ME input field that allows all positive values greater than 0.****

3.3 Example FP scenario with 100-level scale

Suppose we have the following post:

I hate how much I love everybody in the world!!!

Question 1: Please indicate **how hateful** you feel this post is using any number between 0 and 100 where 0 means not hateful and 100 means extremely hateful.

****Show a numerical slider with values between 0 and 100.****

Question 2: SocialNet **removed** the post. The detection system labeled this post as hateful. And so, SocialNet removed the post immediately.

Please indicate whether you agree, disagree, or are neutral about SocialNets decision.

Question 3: Please indicate how much you agree or disagree with SocialNets decision using any positive number from 1 to 100. If you feel neutral about SocialNets decision, select neutral in the field above.

****Show a numerical slider with values between 1 and 100.****

3.4 Example rejection scenario with 100-level scale

Suppose we have the following post:

I hate how much I love everybody in the world!!!

Question 1: Please indicate **how hateful** you feel this post is using any number between 0 and 100 where 0 means not hateful and 100 means extremely hateful.

****Show a numerical slider with values between 0 and 100.****

Question 2: The detection system of SocialNet **was uncertain** about whether the post was hateful or not. An internal moderator needs to look at it. It will remain publicly visible for 24 hours until the moderator decides to remove it when necessary.

Please indicate whether you agree, disagree, or are neutral about SocialNets decision.

Question 3: Please indicate how much you agree or disagree with SocialNets decision using any positive number from 1 to 100. If you feel neutral about SocialNets decision, select neutral in the field above.

****Show a numerical slider with values between 1 and 100.****

4 Analysis

First, we calculate the cost values for the TP, TN, FP, FN, and rejection scenarios in hate speech detection using the survey's results. Second, we analyze whether we can use the results to draw any conclusions by looking at two aspects: reliability and validity.

4.1 Costs

The goal of the complete experiment is to come up with the relative cost values for TP, TN, FP, FN, and rejection scenarios in the context of hate speech detection. The metric from the first part of our research takes these numerical values as its input to calculate the optimal rejection threshold. We do not need to know the absolute cost values but only the relative cost values. For example, if we set all cost values to 1, we retrieve the same optimal rejection threshold as setting all cost values to 1000. Therefore, we need to know the cost ratios between all scenario types. The ME technique provides us with ratio data. All subjects will see the same 40 scenarios. Each scenario simulates a TP, TN, FP, FN, or rejection setting. We use a bipolar scale for question 3 in the survey since we ask the subjects whether they agree, disagree, or are neutral with the decision of SocialNet (tolerating, removing, or rejecting posts). For both scales we will convert disagreement values to negative values, neutral values to 0, and agreement values to positive values. This allows us to calculate the mean value of how much all subjects agree or disagree with the decisions from SocialNet.

For example, to calculate the mean value of all responses to TP scenarios for both scales, we use:

$$\bar{r}_{TP}^{ME} = \frac{1}{n} \sum_{i=1}^n r_i^{ME} \quad \text{where } n \text{ is the total number of ME responses to TP scenarios} \\ \text{and } r_i^{ME} \text{ is the } i\text{th ME response value.}$$
$$\bar{r}_{TP}^{100} = \frac{1}{n} \sum_{i=1}^n r_i^{100} \quad \text{where } n \text{ is the total number of 100-level responses to TP scenarios} \\ \text{and } r_i^{100} \text{ is } i\text{th 100-level response value.}$$

We apply the same calculations for the remaining scenario types. We define the cost values we need for the metric using the mean scores of the TP, TN, FP, FN, and rejection scenarios rated with the ME scale. We will interpret disagreement values as costs and agreement values as gains. We will not use the mean scores of the 100-level scale in our metric since the 100-level scale does not have ratio properties.

4.2 Reliability

Reliability is about whether we can trust our results and if we get consistent results [6]. We do this by mainly looking at the inter-rater reliability. This means that different subjects should give approximately the same

judgements to the same scenarios. We measure the inter-rater reliability using Krippendorff's alpha [12, 10]. We calculate the inter-rater reliability value using the mean response values as calculated in subsection 4.1. We will use the inter-rater reliability scores to compare the ME scale with the 100-level scale. We also study the inter-rater reliability values for the different types of scenarios: TP, TN, FN, FP, and rejection. Other types of reliability, such as test-retest reliability, are not considered in this experiment. Guaranteeing test-retest reliability would require us to redo the complete experiment at a different time for the same subjects. This is infeasible for our Thesis project, given the limited time and budget.

4.3 Validity

Validity is about whether we are measuring the things we want to measure [6]. The main goal of this aspect is to validate if we can use the ME technique to measure subjects' opinions about hate speech detection scenarios. There are multiple types of validity, but we focus mainly on convergent validity (part of construct validity), content validity, and face validity [6]. Construct validity checks whether there is an agreement between a theory and a measurement device or procedure [6]. Convergent validity is about the correlation between different types of measures to see if they measure the same phenomenon [6]. Content validity is about letting external experts review the proposed research questions and procedure [6]. Face validity is the most intuitive and subjective type of validity and is about if and why we think the questions and proposed procedures are valid [6].

We analyze convergent validity by comparing the mean scores from subsection 4.1 between the two scales. We can verify that they measure the same phenomenon by analyzing the correlation between the scales. However, we can expect a low correlation since the ME scale is a (normalized) unbounded scale and the 100-level scale is bounded. Nevertheless, we think that both scales will give similar results, meaning that high ME responses should correspond to high 100-level scale responses and low ME responses to low 100-level scale responses. To guarantee content validity, we let external experts (the supervisors of this Thesis project) check this pre-registration report. We guarantee face validity by discussing whether the ME technique gives us the expected results. We exclude other forms of validity from this experiment because they either are irrelevant or infeasible. For example, external validity is about the degree to which the findings can be generalized to other settings or groups. We would have to experiment with multiple groups with different demographic and age characteristics to guarantee external validity. We left this for future work to investigate since this is out of scope for our Thesis project. Despite that, we think that people of different ages and demographic characteristics perceive hate differently since people have other norms and values in various parts of the world. We believe that if we conduct this experiment using different groups of subjects, then we might retrieve different cost values. Therefore, we decided not to create too many participant inclusion criteria but take a random sample of global social media users.

A Presentation texts

A.1 Consent

You are being invited to participate in a research study titled "Costs of predictions in hate speech detection". This study is being done by Philippe Lammerts from the TU Delft.

The purpose of this research study is to find out what social media users think of different scenarios of hate speech detection on social media. It will take you approximately 25 minutes to complete. These scenarios consist of two things. First, we show a specific social media post that can be either hateful or not hateful. You need to indicate how hateful you feel that the post is. Second, we explain how the social media platform dealt with this post. You need to indicate whether you agree/disagree/are neutral about the platform's decision. The results of the survey will be used in my thesis.

As with any online activity, the risk of a breach is always possible. To the best of our ability, your answers in this study will remain confidential. We will minimize any risks by making this survey completely anonymous. Therefore, please do not provide any personal information anywhere. The anonymous results might be shared publicly in the future.

Your participation in this study is entirely voluntary, and you can withdraw at any time.

Feel free to contact me with any questions or feedback you might have: p.m.lammerts@student.tudelft.nl

If you understand and agree with the above information and consent to take part in this study, then you can check the checkbox and click on the 'Next' button to start the survey.

A.2 Short introduction ME-100

- This survey consists of two parts.
- In each part, you will be presented with a series of different scenarios.
- For each scenario, you need to answer two questions using a specific scale.
- We will explain the exact instructions later.
- But first, we will let you familiarize yourself with a scale called Magnitude Estimation.

A.3 Short introduction 100-ME

- You will be presented with a series of different scenarios.
- Each scenario describes a situation of a social media user who wants to post a specific message on a fictional social media platform we now call SocialNet.
- These posts can be neutral or contain hateful content (sexist or racist posts).
- SocialNet uses automated detection systems for detecting hate speech.

A.4 Introduction

You will be presented with a series of different scenarios.

- Each scenario describes a situation of a social media user who wants to post a specific message on a fictional social media platform we now call SocialNet.
- These posts can be neutral or contain hateful content (sexist or racist posts).
- SocialNet uses automated detection systems for detecting hate speech.
- SocialNet can do one of the following things for each scenario:
 1. SocialNet **removed** the post. The detection system labeled this post as hateful. And so, SocialNet removed the post immediately.
 2. SocialNet **tolerated** the post. The detection system labeled this post as not hateful. And so, it remains publicly visible.
 3. The detection system of SocialNet **was uncertain** about whether the post was hateful or not. An internal moderator needs to look at it. It will remain publicly visible for 24 hours until the moderator decides to remove it when necessary.

For each scenario, you need to answer two questions:

A.5 100-level scale explanation

- First, your task is to indicate how hateful you feel that a post is.
 - You do this by assigning any number between 0 and 100, where 0 means 'not hateful' and 100 means 'extremely hateful'.
 - Try to make each number match the intensity as you perceive it.
 - Please try to not only use appropriate numbers but also avoid restricting your choice of numbers from 1 to 10. ([3] recommended to include this instruction)
- Second, your task is to tell how much you agree or disagree with the decision from SocialNet.
 - If you feel neutral about SocialNets decision, this value will be equal to 0.
 - If you (dis)agree with the decision, you need to indicate how much you (dis)agree with it by assigning any number between 1 and 100.
 - A large number means you (dis)agree with it a lot, while a small number means you (dis)agree with it a little.
 - Try to make each number match the intensity as you perceive it.

- Please try to not only use appropriate numbers but also avoid restricting your choice of numbers from 1 to 10. ([3] recommended to include this instruction)

A.6 ME scale explanation (inspired by [14])

- First, your task is to indicate how hateful you feel that a post is by assigning any number greater or equal to 0, where 0 means that a post is not hateful at all.
 - Assign any number that seems appropriate to you.
 - If you feel that the current post is twice as hateful as the previous one, you need to assign a number that is twice as large as the previous number.
 - Or, if you feel that the current post is half as hateful as the previous one, you need to assign a number that is half as large as the previous number.
 - If you feel that something is not hateful at all, you need to assign a 0.
- Second, your task is to tell how much you agree or disagree with the decision from SocialNet.
 - If you feel neutral about SocialNets decision, this value will be equal to 0.
 - If you (dis)agree with the decision from SocialNet, you need to assign any number that is greater or equal to 1 that reflects how much you (dis)agree with the decision.
 - Assign any number that seems appropriate to you.
 - If you (dis)agree twice as much with the current decision as with the previous one, you need to assign a number that is twice as large as the previous number.
 - Or, if you (dis)agree half as much with the current decision as with the previous one, you need to assign a number that is half as large as the previous number.

You can use any number or decimal you want for both tasks, but make each assignment proportional to your subjective impression.

A.7 Attention check

This question is an attention check. You must select 'Blue' here.⁵

A.8 Training phase ME

As a warm-up task, to familiarize you with magnitude estimation, you will be shown a sequence of five lines, one at a time.

- For each line, enter a number into the text box below the displayed line. This number should reflect your perception of the length of the line. You may use any numbers that seem appropriate to you whole numbers or decimals. However, you may not use negative numbers or zero.
- For each subsequent line, enter a number that reflects your perception of its length, relative to the previous line. For example, if you feel that the current line is twice as long as the previous, then you should assign a number that is twice as large as the number you used previously.

Dont worry about running out of numbers there will always be a larger number than the largest you use, and a smaller number than the smallest you use. Note: The magnitude estimation scores are **not** intended to be an estimate of the length in any particular measurement units, such as centimeters. [12]

A.9 ME calibration

You are now finished with the scenarios using the magnitude estimation scale.

In the **first question** of each scenario, you had to indicate how hateful you felt that a post was.

We now ask you to estimate **what numbers correspond to what labels**.

⁵Based on <https://researcher-help.prolific.co/hc/en-gb/articles/360009223553>

- Try to recall which numbers you used before when answering the questions and assign each label a number that reflected your subjective feeling.
- For example, if you found a post during the scenarios 'extremely hateful' and rated it using the number 300, then assign 300 to the label 'extremely hateful'.
- For example, if you found a post during the scenarios 'slightly hateful' and rated it using the number 8, then assign 8 to the label 'slightly hateful'.

****Show 'Not hateful (equals 0)', 'Slightly hateful', 'Hateful', and 'Extremely hateful' labels with ME inputs****

In the **second question** of each scenario, you had to indicate how much you agreed or disagreed with how SocialNet dealt with the social media post.

We now ask you to estimate **what numbers correspond to what labels**.

- Try to recall which numbers you used before when answering the questions and assign each label a number that reflected your subjective feeling.
- Use negative values for the 'strongly disagree', 'disagree', and 'somewhat disagree' fields.
- Use positive values for the 'strongly agree', 'agree', and 'somewhat agree' fields.
- For example, if you strongly agreed with a decision during the scenarios and rated it using the number 300, then assign 300 to the label 'strongly agree'.
- For example, if you somewhat disagreed with a decision during the scenarios and rated it using the number 8, then assign -8 to the label 'somewhat disagree'.

****Show 'Strongly disagree', 'Disagree', 'Somewhat disagree', 'Neutral (equals 0)', 'Somewhat agree', 'Agree', and 'Strongly agree' labels with ME inputs****

References

- [1] I. E. Allen and C. A. Seaman. Likert scales and data analyses. *Quality progress*, 40(7):64–65, 2007.
- [2] A. Balayn, J. Yang, Z. Szlavik, and A. Bozzon. Automatic identification of harmful, aggressive, abusive, and offensive language on the web: A survey of technical biases informed by psychology literature. *ACM Transactions on Social Computing (TSC)*, 4(3):1–56, 2021.
- [3] E. G. Bard, D. Robertson, and A. Sorace. Magnitude estimation of linguistic acceptability. *Language*, pages 32–68, 1996.
- [4] V. Basile, C. Bosco, E. Fersini, N. Debara, V. Patti, F. M. R. Pardo, P. Rosso, M. Sanguinetti, et al. Semeval-2019 task 5: Multilingual detection of hate speech against immigrants and women in twitter. In *13th International Workshop on Semantic Evaluation*, pages 54–63. Association for Computational Linguistics, 2019.
- [5] H. N. Boone and D. A. Boone. Analyzing likert data. *Journal of extension*, 50(2):1–5, 2012.
- [6] K. Fitzner. Reliability and validity a quick review. *The Diabetes Educator*, 33(5):775–780, 2007.
- [7] M. W. T. H. D. Gold and T. Zesch. Do women perceive hate differently: Examining the relationship between hate speech, gender, and agreement judgments. 2018.
- [8] T. Gröndahl, L. Pajola, M. Juuti, M. Conti, and N. Asokan. All you need is "love" evading hate speech detection. In *Proceedings of the 11th ACM workshop on artificial intelligence and security*, pages 2–12, 2018.
- [9] K. Hendrickx, L. Perini, D. Van der Plas, W. Meert, and J. Davis. Machine learning with a reject option: A survey. *arXiv preprint arXiv:2107.11277*, 2021.
- [10] K. Krippendorff. Reliability in content analysis: Some common misconceptions and recommendations. *Human communication research*, 30(3):411–433, 2004.
- [11] M. Lodge and B. Tursky. Comparisons between category and magnitude scaling of political opinion employing src/cps items. *American Political Science Review*, 73(1):50–66, 1979.

- [12] E. Maddalena, S. Mizzaro, F. Scholer, and A. Turpin. On crowdsourcing relevance magnitudes for information retrieval evaluation. *ACM Transactions on Information Systems (TOIS)*, 35(3):1–32, 2017.
- [13] M. McGee. Master usability scaling: magnitude estimation and master scaling applied to usability measurement. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 335–342, 2004.
- [14] H. R. Moskowitz. Magnitude estimation: notes on what, how, when, and why to use it. *Journal of Food Quality*, 1(3):195–227, 1977.
- [15] J. Murray. Likert data: what to use, parametric or non-parametric? *International Journal of Business and Social Science*, 4(11), 2013.
- [16] G. Norman. Likert scales, levels of measurement and the laws of statistics. *Advances in health sciences education*, 15(5):625–632, 2010.
- [17] J. S. Olson and W. A. Kellogg. *Ways of Knowing in HCI*, volume 2. Springer, 2014.
- [18] K. Roitero, E. Maddalena, G. Demartini, and S. Mizzaro. On fine-grained relevance scales. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 675–684, 2018.
- [19] S. S. Stevens. The direct estimation of sensory magnitudes: Loudness. *The American journal of psychology*, 69(1):1–25, 1956.
- [20] A. E. Van’t Veer and R. Giner-Sorolla. Pre-registration in social psychologya discussion and suggested template. *Journal of experimental social psychology*, 67:2–12, 2016.
- [21] Z. Waseem and D. Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop*, pages 88–93, 2016.