

“...I can show what I really like.”: How Quadratic Voting better align true preferences than Likert Scale Surveys

A clear and well-documented \LaTeX document is presented as an article formatted for publication by ACM in a conference proceedings or journal publication. Based on the “acmart” document class, this article presents and explains many of the common variations, as well as many of the formatting elements an author may use in the preparation of the documentation of their work.

CCS Concepts: • **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

Additional Key Words and Phrases: datasets, neural networks, gaze detection, text tagging

ACM Reference Format:

. 2020. “...I can show what I really like.”: How Quadratic Voting better align true preferences than Likert Scale Surveys. In *CSCW ’20: The 23rd ACM Conference on Computer-Supported Cooperative Work and Social Computing*, Oct 17 – 21, 2020, Virtual. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Likert scale survey is one of the most widely used methods to obtain participant’s opinion in the realm of human-computer interaction. Survey participants would express a rating across series of measurements — *Very agree to very disagree* or *On a scale of 1 to 5* — for a listed statement. Very often, these opinions help researchers or decision makers uncover consenses across a group of people.

However, there had been finding of how researchers can easily misuse likert scale surveys either applying incorrect analysis methods [2] or misinterpreting the analysis results [6, 13] leading to questionable findings. In addition, many research papers do not explain the rational behind the use of Likert scale surveys. In a community that adopted Likert scale surveys almost as the defacto standard, we ask a fundamental question: “Is Likert-scale survey the ideal method to measure collective attitudes for decision making?”

We begin by exploring one type of question in collective decision making that aims To elicit user preferences among K options. Research agencies, industry labs or independent researchers often want understand how to better allocate resources. For example, ordinal scale polls were designed to understand public opinions on government policy [1] because there are limited funding. Companies deploy online surveys to understand how product users feel about the features and services that needs further improvements because companies has limited time to develop the next release. Physical surveys can be found in shopping centers to collect individual’s experiences for products on the shelf because there are limited shelves. All these example demonstrated how surveys are often tied to making decisions by gathering consensus from surveying individual’s attitudes.

In this study, we look at an alternative method called Quadratic Voting (QV). Published in 2015, Weyl at al. [14] proposed Quadratic voting as a voting mechanism with approximate Pareto efficiency.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CSCW ’20, Oct 17 – 21, 2020, Virtual

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-9999-9/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

Under this voting mechanism, voters were initially given a fixed amount of voice credits (VC). With the credits, individuals can purchase any number of votes to support any of the statements listed on the ballot. However, the cost of each vote increases quadratically when voted toward the same option. The authors proved that this mechanism is more efficient at making a collective decision because it minimizes welfare loss. Since 2015, a few studies compared Likert scaled surveys with QV empirically and theoretically [10, 15]. Cavaille et. al argues that QV outperforms Likert-scale surveys among a set of political and economic issues [3]. Despite these findings, we are not aware of related works that compare Likert scale surveys and QV with participants' underlying true preferences. Therefore, it is unclear whether or not and in what degree does QV results align with participants behaviors. In addition, no current work, to the best of our knowledge, deployed QV in the area of HCI.

To be more specific, we ask the following research questions:

RQ1. How does results from QV, Likert scaled survey align with people's behavior when surveying societal issues?

RQ2. How does results from QV, Likert-scale align with people's behavior when placed in an HCI context?

RQ3. How does different number of voice credits impact results of QV empirically?

RQ4. What are some qualitative insights that can be observed when participants vote under QV?

To answer these research questions, we designed two experiments. The first experiment, designed to answer RQ1 and RQ3, is a between-subject study where participants express their attitudes among a set of societal causes using QV and Likert-scaled surveys and then donate to organizations relevant to these organizations. The second experiment created an HCI study environment, aimed to answer RQ2, where participants were asked about opinions among different video elements and their opinions using QV and Likert-scaled surveys. Our results showed that both experiments support QV in providing a clean and efficient way compared to Likert scale surveys at eliciting participant's true preferences.

Contributions Our work made several contributions to the research community. First, we proved empirically the use of QV outperforms Likert scale survey when conducting "choosing one in K" experiments. Second, we showed that the usability of QV is transferrable from a general domain to HCI. Third, we designed a bayesian model that facilitates the comparison of Likert scale surveys, QV and behaviors. Fourth, we designed an online experiment to mimic real life HCI-related decision making. And finally, we provided the source code of our easy to deploy, interactive web platform for QV to the community.

2 RELATED WORKS

In this section, we laid out the related works for QV and Likert scaled surveys.

2.1 QV

- qv theory
- qv in practice

2.2 Likert

- development in likert
- application of likert in HCI
- challenges in likert

3.2.1 *Selection of the societal causes.*

3.2.2 *System Architecture and Interface.* The voting system is constructed using Python Flask for the back-end, Angular for front-end and MongoDB for database storage. The experiment source code is publicly available ¹ and the QV interface is also provided as a stand-alone repository ². In this subsection, we focus on the QV interface.

The QV interface, shown in graph Y, consists of three major sections. The first section contains definitions of QV and the prompt of the task. The second section shows a list of option with a plus and minus button to its left. Buttons are disabled if the number of voice credit does not permit the next vote. A bar on the right of the option shows the proportion of voice credits used to that option. The final section is a floating summary at sticks to the bottom of the page. It contains a visualization of the total number of credits and the remaining credits.

3.3 Experiment 2

The second experiment extends upon the first one, in which it examines whether Quadratic Voting betters at aligning people's actual preferences compared to a Likert-scaled survey in an HCI setting. Different from political and public-opinion surveys, testing participants' preference in interface design and user experience is much more non-trivial. Thus we developed a buy-back mechanism and observe participants' behaviors as their true preference. This experiment also acts as a concrete example as to how QV can be incorporated in HCI.

3.3.1 *Choice of HCI Research Question.* Research on video and audio quality from the lens of HCI has been a relatively mature. Contributions has been made to fields like multi-media conferencing [17], video-audio perception [4, 9] and more specifically trade-offs between video and audio elements under network monetary constraints [8, 12].

Oeldorf-Hirsch et al. [12] conducted a study, covering the widest range of elements to the best of our knowledge, to understand how users with bandwidth constraints made trade-offs between video and audio elements. They examined participants' attitude between three video bit rates, three video frame rates and two audio sampling rates across three types of video content. Participants were asked to rate the overall quality, video quality, audio quality and enjoyment level on a 5-point Likert scale in each condition. Conclusion were drawn using mean and standard deviation of the survey results. This is a typical study where the goal is to find 1 or some of the K elements to choose from when under constraint. In our second experiment, we expand this study to collect people's preference among a wider range of video and audio elements and compare how Likert-scaled survey and QV reflects people's true perception preferences.

3.3.2 *Experiment 2 Design.* In our experiment, we included a total of five video and audio element that will impact a video. These elements include video and audio package loss rate, determining whether the audio or video stutters; video resolution and audio sampling rate effecting the quality of video and audio; and video-audio synchronization. We selected a few segment of weather broadcasting from a news channel as the content of our video. Weather broadcasts usually convey information via both visual and audio channels, appeal to a wide array of audiences, and do not require prior knowledge to understand.

To ensure the ecological validity of the experiment, we situated the comparison of different video and audio elements in a hypothetical scenario in which the participant is a manager of a weather reporting news station. As the manager, the participant was asked to rate the importance of each

¹Not yet public

²<https://github.com/hank0982/QV-app>

video and audio elements with the goal to maximize customer understanding of the context where network is of low bandwidth and that the weather broadcast cannot be shown in its best quality.

We designed a between-subject study with three groups of 60 participants. After the participants agreed with the consent form, all three group of participant were presented an example weather broadcast segment with controls of the five video and audio elements under the video shown in figure M. All five elements were set to sub-optimal by default, making the content near incomprehensible. Participants can alternate the five elements in any combination, to see how elements impact to the video.

Once participants think that they had a grasp of how different elements impact a video, the first group of participants then completed a 5-point Likert-scaled survey while the second group of participants completed a QV survey with K voice credits ³, asking their opinions on the importance of the 5 video and audio elements in a weather broadcast under a low bandwidth environment. The third group of participants were asked to perform a buy-back task for a bad-quality advertising video.

The buy-back task mimics a rational customer's behavior: buying essential tools to complete some given task. Participants were told that as the manager of the weather broadcast agency, they need to verify if their viewer can understand the content of the video. Therefore, the goal of their task is to correctly answer a set of multiple choice questions to make sure that they correctly comprehend the video. Given the video with sub-optimal video, participants were given a budget of \$30 to purchase some or all of the features back. To ensure incentive-compatibility of the participants' buy-back actions, we offered to pay the participants their own remaining amount from the \$30 budget through a lottery under the condition that they answered 80% of the multiple choice questions correctly, [missing probability] version of a new weather broadcast video adjusted by their buy-back choices. These questions contained factual questions such as, "What is the weather of Chicago?", "What is the highs and lows of San Diego", "Which of the follow cities got colder?". Participants were shown three example questions before the buy-back task to assist their decision. Participants can replay the video with their adjustments while answering the questions to ensure that participants do not require memorization. There will be a 5 minute timer to minimize the impact of replaying the video. With this design, participants would try their best to make the video comprehensible based on their opinions on which feature(s) was most needed at the lowest cost.

In the given weather broadcast video, there were 4 levels of quality for each of the 5 elements. By default, the video set to the lowest level for all elements before any adjustment occurred.

- (1) Audio Package Loss Repaired with Silence (package loss rate) [17]: 20%, 10%, 5%, 0%
- (2) Video Package Loss (package loss rate) [5]: 20%, 8%, 4%, 0% (20, 8.3, 3.3, 0)
- (3) Audio Sampling Rate [11, 12]: 8kHz, 11kHz, 16kHz, 48kHz
- (4) Video Resolution [7, 12]: 120x90, 168x126, 208x156, 240x180
- (5) Video-audio Synchronization (time video behind audio) [16]: 240ms, 200ms, 160ms, 0ms (new: 1850, 1615, 1050, 0)

In the buy-back task, each level of improvement for one feature costs \$2. It would cost the entire budget of \$30 to buy all levels of every feature back. Hence, the option of buying back everything was given to the participants, in return, there would be no extra payoff remaining for the participant.

Similar to the first experiment, the money spent on each feature during the buy-back task are considered as the true preference the population had towards the 5 video and audio features. The results from the Likert-scaled surveys and QV survey were then compared to the population's true preference to see how different they were.

³K is decided from experiment 1

4 RESULTS

4.1 Experiment 1

4.1.1 report results...

REFERENCES

- [1] 2018. 2018 Midterm Voters: Issues and Political Values. *Pew research center* (2018).
- [2] Phillip A Bishop and Robert L Herron. 2015. Use and misuse of the Likert item responses and other ordinal measures. *International journal of exercise science* 8, 3 (2015), 297.
- [3] Charlotte Cavaillé, Daniel L Chen, and Karine Van Der Straeten. 2018. Towards a General Theory of Survey Response: Likert Scales Vs. Quadratic Voting for Attitudinal Research. (2018).
- [4] Sherry Y Chen, Gheorghita Ghinea, and Robert D Macredie. 2006. A cognitive approach to user perception of multimedia quality: An empirical investigation. *International Journal of Human-Computer Studies* 64, 12 (2006), 1200–1213.
- [5] Mark Claypool and Jonathan Tanner. 1999. The effects of jitter on the peceptual quality of video. In *Proceedings of the seventh ACM international conference on Multimedia (Part 2)*. ACM, 115–118.
- [6] Susan Jamieson et al. [n.d.]. Likert scales: how to (ab) use them. *Medical education* 38, 12 ([n. d.]), 1217–1218.
- [7] Hendrik Knoche, John D Mccarthy, and M Angela Sasse. 2008. How low can you go? The effect of low resolutions on shot types in mobile TV. *Multimedia Tools and Applications* 36, 1-2 (2008), 145–166.
- [8] Andreea Molnar and Cristina Hava Muntean. 2013. COMEDY: Viewer trade off between multimedia quality and monetary benefits. In *2013 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 1–7.
- [9] Andreea Molnar and Cristina Hava Muntean. 2015. Assessing learning achievements when reducing mobile video quality. *Journal of Universal Computer Science* 21, 7 (2015), 959–975.
- [10] Ryan Naylor et al. 2017. First year student conceptions of success: What really matters? *Student Success* 8, 2 (2017), 9–19.
- [11] Peter Noll. 1993. Wideband speech and audio coding. *IEEE Communications Magazine* 31, 11 (1993), 34–44.
- [12] Anne Oeldorf-Hirsch, Jonathan Donner, and Edward Cutrell. 2012. How bad is good enough?: exploring mobile video quality trade-offs for bandwidth-constrained consumers. In *Proceedings of the 7th Nordic Conference on Human-Computer Interaction: Making Sense Through Design*. ACM, 49–58.
- [13] Godfrey Pell. 2005. Use and misuse of Likert scales. (2005).
- [14] Eric A Posner and E Glen Weyl. 2018. *Radical markets: Uprooting capitalism and democracy for a just society*. Princeton University Press.
- [15] David Quarfoot, Douglas von Kohorn, Kevin Slavin, Rory Sutherland, David Goldstein, and Ellen Konar. 2017. Quadratic voting in the wild: real people, real votes. *Public Choice* 172, 1-2 (2017), 283–303.
- [16] Ralf Steinmetz. 1996. Human perception of jitter and media synchronization. *IEEE Journal on selected Areas in Communications* 14, 1 (1996), 61–72.
- [17] Anna Watson and Martina Angela Sasse. 1996. Evaluating audio and video quality in low-cost multimedia conferencing systems. *Interacting with Computers* 8, 3 (1996), 255–275.