

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

A clear and well-documented L^AT_EX 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, 5 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Likert scale survey is one of the most widely used methods to obtain the participant’s opinion in the realm of human-computer interaction. Survey participants would express a rating across a 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 findings of how researchers can easily misuse Likert scale surveys either applying incorrect analysis methods [3] or misinterpreting the analysis results [10, 17] 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 to understand how to better allocate resources. For example, ordinal scale polls were designed to understand public opinions on government policy [1] because there is limited funding. Companies deploy online surveys to understand how product users feel about the features and services that needs further improvements because companies have limited time to develop the next release. Physical surveys can be found in shopping centers to collect an individual’s experiences for products on the shelf because there are limited shelves. All these examples demonstrated how surveys are often tied to making decisions by gathering consensus from surveying individual’s attitudes.

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>

In this study, we look at an alternative method called Quadratic Voting (QV). Published in 2015, Weyl et al. [18] proposed Quadratic voting as a voting mechanism with approximate Pareto efficiency. 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 [14, 19]. Cavaille et al. argues that QV outperforms Likert-scale surveys among a set of political and economic issues [4]. 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 do different amounts 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 generic domain to HCI. Third, we designed a bayesian model that facilitates the comparison of Likert scale surveys, QV, and behaviors. Fourth, we developed 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.

Design Implication TODO. Talk about interface, future work and insights.

2 RELATED WORKS

In this section, we first explain the challenges that Likert faces and then describe related works of QV.

2.1 Likert Scale Surveys

Likert-scale survey, is a commonly used method to collect participant's opinions because of its ease of use. These surveys are deployed to validate findings or clarify hypotheses [12, 15] in HCI. They are also used to verify or uncover the user's needs. The original Likert scale surveys were invented by Rensis Likert in 1932, which utilizes step-intervals from one attitude to the next on the scale [13]. Researchers today design 3, 5, 7, or even 12-point Likert scale evaluation surveys to accommodate

different uses [7, 8]. In addition, these surveys can also use verbal descriptions to demonstrate an ordinal scale. Some researchers even developed alternative forms of Likert scale such as slider scales [20] or phrase completions [9], which aimed to circumvent some of the shortcomings of the traditional Likert scale.

There are, however, widespread controversies in the community on when, why, and how to use it [3]. For instance, researchers can misuse statistical methods, such as using mean and standard deviations [10], to understand outcomes when working with an ordinal metric. Quantifying aggregated Likert scale surveys, such as explaining what “agree and a half” means can also be unclear. Besides, as options on the survey can be stepwise, one should not assume scales to be equally divided, which can be confusing. In other words, strongly agree and agree can be different compared to neutral and agree [5, 10].

An empirical study by Quarfoot et al. identified another challenge where people exaggerate their views when filling out political surveys [19]. In this study, participants often express strong polarized opinions or have no opinion at all, making it hard to form optimal conclusions [18]. This occurrence was theoretically proved by Cavaille et al. [4], where respondents tend to overstate their values if they want to influence the results through the survey. These challenges motivated us to understand whether QV can fill the gap and provide a more accurate measurement for collective decision making.

2.2 Quadratic Voting

Quadratic voting originated with the argument that current one-person-one-vote system can easily bias toward the majority’s opinion and omitting the minority’s votes [18]. This phenomenon is termed as the tyranny of the majority where the democratic decision does not take care of those in need. In other words, these types of voting does not allow fine-grain responses to the options they were to vote [16]. Some voting mechanisms tried to resolve this by introducing rank-based voting in which voters would decide how they rank the options while submitting their opinions. This mechanism can however suffer from Condorcet’s paradox where results can be suboptimal because the ranks of the voters might not be transitive [16]. Many other voting mechanisms suffer from similar issues.

This triggers the development of Quadratic Voting created by Weyl et al. to overcome traditional voting challenges [18]. QV tries to capture a cost for the voter when he or she made a particular decision by voting toward specific options. This “price-taking” equilibrium helps participants maximize their utility using their votes. This is theoretically proven by Lalley et al. [11] and showed that there exists an approximate structure of Bayes-Nash equilibria.

In order for QV to capture the voice of the minority and allocate the correct cost to the votes QV has the following mechanism: Consider collecting N participants voting, each participant is entitled to K voice credits. Participants can express their binary opinion (for or against) each option o_i within a set of options O listed on the ballot. Participants can purchase any number of votes $v_i \in \mathbb{R}$ vetted toward any of the options o_i . However, to vote v_i votes toward o_i , participants have to spend v_i^2 voice credits, billed toward their k credits. The outcome of QV would be the ranks of the sum among the total votes for any option $\sum V_{oi}$ across all N participants.

To use an analogy, Suppose every voter has a bank account with a fixed amount of money, say 100 dollars. On a ballot, there are ten statements. Voters can now buy votes using their money in the account. However, for each statement, the cost of the votes increases quadratically. For example, voting two for on the first option would cost the voter four votes; voting three against on the fifth option costs the voter nine votes, and so on. This means that the more votes one devote to an option, the more costly it is to do so, forgoing the opportunity the voter had to vote for other options.

2.3 QV in the wild

After QV was proposed, Quarfoot et al. conducted an empirical study to understand how QV results compare to Likert scale surveys. They recruited 4500 participants to survey an individual's opinions across ten public policy using either Likert-style questionnaire, QV survey, or both. The study found that the number of people who voted for the same number of votes for the options distributed normally and consistently across all options. This differs from results from Likert scaled surveys completed from the same group of participants, where results are either heavily skewed or polarized "W-shaped" distribution. Researchers also saw individuals spent more time expressing their opinion and reveal a more fine-grain attitude toward the policies. Thus, the study concludes that QV provides a clearer picture of public opinion to policymakers [19].

The work by Quarfoot et al., however, only used mean and z-scores, to compare the final aggregated results across the two methods. In addition, the design on the policies have a strong tendency for voters to agree or disagree on the extreme, such as one's opinion on "same-sex marriage". Thus, little do we know if QV produces different results than Likert scale surveys if the options are less competitive, for example, choosing one's favorite ice cream flavor.

Another empirical study was applied to education by Ryan Naylor [14]. The author again used QV and Likert to understand essential elements among a list of factor that impacts students' success in universities. Results showed that QV provided more insights, such as distinguishing good-to-have factors from must-have elements. These factors are not heated debated controversies compared to public policies in the previous studies and they are independent elements that do not require students to make trade-offs. For example, students can have a sense of "belonging" and a sense of "achievement" at the same time.

To the best of our knowledge, there does not exist an empirical study that focused on investigating how QV and Likert perform under the condition of selecting one in K options. This setting was recently discussed in a theoretical work by Eguia et al. [6], who claims that QV is still in favor of resolving budget-constrained for risk natural agents to figure out an efficient decision across multiple alternatives as a collective choice problem. We aim to complete this missing piece of the puzzle. Further, we are not aware of any work that studies alignment between participants' actual beliefs and QV surveys. Existing research pointed out possible fallacy exists with self-reporting [2, 21]. Thus, we aim to understand how QV and Likert scale surveys align with the agent's true beliefs. We also want to test whether the total number of voice credits impacted the results of QV. Finally, we believe that no HCI research utilized QV to form design decisions.

3 EXPERIMENT 1

3.1 Methodology

3.2 Experiment 1 Results

4 EXPERIMENT 2

4.1 Methodology

4.2 Experiment 2 Results

5 DISCUSSION

6 CONCLUSION

REFERENCES

- [1] 2018. 2018 Midterm Voters: Issues and Political Values. *Pew research center* (2018).
- [2] Theo Araujo, Anke Wonneberger, Peter Neijens, and Claes de Vreese. 2017. How much time do you spend online? Understanding and improving the accuracy of self-reported measures of Internet use. *Communication Methods and Measures* 11, 3 (2017), 173–190.

- [3] 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.
- [4] 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).
- [5] Diane R Edmondson. 2005. Likert scales: A history. In *Proceedings of the 12th conference on historical analysis and research in marketing (CHARM)*. 127–133.
- [6] Jon X Eguia, Nicole Immorlica, Katrina Ligett, E Glen Weyl, and Dimitrios Xefteris. 2019. Quadratic voting with multiple alternatives. *Available at SSRN 3319508* (2019).
- [7] Kraig Finstad. 2010. Response interpolation and scale sensitivity: Evidence against 5-point scales. *Journal of Usability Studies* 5, 3 (2010), 104–110.
- [8] Kate J Garland and Jan M Noyes. 2008. Computer attitude scales: How relevant today? *Computers in Human Behavior* 24, 2 (2008), 563–575.
- [9] David R Hodge and David Gillespie. 2003. Phrase completions: An alternative to likert scales.(Note on Research Methodology). *Social Work Research* 27, 1 (2003), 45–56.
- [10] Susan Jamieson et al. [n.d.]. Likert scales: how to (ab) use them. *Medical education* 38, 12 ([n. d.]), 1217–1218.
- [11] Steven P Lalley and E Glen Weyl. 2018. Quadratic voting: How mechanism design can radicalize democracy. In *AEA Papers and Proceedings*, Vol. 108. 33–37.
- [12] David Ledo, Steven Houben, Jo Vermeulen, Nicolai Marquardt, Lora Oehlberg, and Saul Greenberg. 2018. Evaluation strategies for HCI toolkit research. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [13] Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of psychology* (1932).
- [14] Ryan Naylor et al. 2017. First year student conceptions of success: What really matters? *Student Success* 8, 2 (2017), 9–19.
- [15] A Ant Ozok. 2009. Survey design and implementation in HCI. *Human-Computer Interaction: Development Process* 253 (2009).
- [16] Eric Pacuit. 2019. Voting Methods. In *The Stanford Encyclopedia of Philosophy* (fall 2019 ed.), Edward N. Zalta (Ed.). Metaphysics Research Lab, Stanford University.
- [17] Godfrey Pell. 2005. Use and misuse of Likert scales. (2005).
- [18] Eric A Posner and E Glen Weyl. 2018. *Radical markets: Uprooting capitalism and democracy for a just society*. Princeton University Press.
- [19] 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.
- [20] Catherine A Roster, Lorenzo Lucianetti, and Gerald Albaum. 2015. Exploring slider vs. categorical response formats in web-based surveys. *Journal of Research Practice* 11, 1 (2015), 1.
- [21] Lynn Vavreck et al. 2007. The exaggerated effects of advertising on turnout: The dangers of self-reports. *Quarterly Journal of Political Science* 2, 4 (2007), 325–343.