

1 Author Response: Explainable Voting #8353

2 Review 1

3 We discuss the impact of our results for machine learning applications in lines 28–48 of the submission. The reflections
4 by Reviewer 4 in points 2 and 8 are a good expression of how we hope to enrich techniques for explainable ML.

5 Review 2

6 A major question raised by your review is: how much flexibility do we have in choosing our axiom systems? The
7 answer is that we are heavily constrained. It is possible to write down trivial axiom systems such as your INIT' example,
8 but of course these axioms are not normatively appealing. For other axiom systems, it is necessary to establish that they
9 *uniquely characterize* the voting rule in question, since otherwise there will be instances where the algorithm fails to
10 find an explanation. Very few characterization results are known for voting, even for rules other than Borda. (Worse,
11 most known characterizations rely on limit arguments, which does not lead to explanations in our sense.) Arguably, in
12 our paper, we have captured most of the important ones in a single framework.

13 Given this background, our choice of axioms for the Borda rule is not arbitrary: the axiom system we use (up to minor
14 variations) is the only one known to characterize Borda (without limit arguments), and luckily the axioms used have
15 considerable normative appeal. While, as you note, the INIT axioms could be seen as arbitrary, in the Borda case, they
16 seem well-motivated via simple symmetry arguments.

17 Regarding using the simple embedding of Borda into space \mathbb{N}^m to get shorter Borda explanations: indeed this would be
18 possible, but in our framework the EMB axiom would then require that $f(R_1) = f(R_2)$ whenever profiles R_1 and R_2
19 have the same Borda scores, and this seems too specific. While we don't think this specific embedding will lead to
20 convincing explanations, the thought process you engaged in as a reader is an example of a hope we express in the
21 discussion: "The good news is that Theorem 2 can help identify new axiomatizations that lead to short explanations." In
22 other words, the possibility of new embeddings that would lead to new characterizations and simpler explanations is a
23 clear strength of our general framework, not a weakness.

24 Regarding the definition of "asymptotically weaker": you are correct that there is a natural version of this definition
25 for general axiom systems. However our definition specific to \mathcal{S}_{emb} exploits some additional freedom we have in our
26 specific set-up (namely, that we are allowed to use an unlimited number of ADD and EMB axioms), and we need this
27 freedom in our application of the framework to Borda.

28 Regarding "trade-off theorems" and measuring the "strength" of classes of axiom systems: this is an intriguing idea, but
29 it is not clear how to formalize the notion of a "class" of axiom systems. For the voting context, we are again limited by
30 the small number of axiom systems (with normative appeal) that are known to characterize common voting rules.

31 Review 3

32 You mention that we do not include empirical examples. For illustration purposes, we did include (in the supplementary
33 material) a sample explanation of Borda applied to the mayoral election in Burlington, VA. Since by their nature the
34 generated explanations follow a common pattern, we did not think it would be instructive to give many more examples.
35 Quantitatively speaking, we did not see sufficiently promising avenues for empirical exploration: since our bounds
36 provably apply to almost all instances, an empirical evaluation will not reveal that shorter explanations are possible in
37 practice. [And this prediction was confirmed by preliminary experiments we have run.]

38 Review 4

39 We were amazed by your in-depth review, full of great suggestions. Thanks for taking the time to think deeply about
40 our paper. Responding (much too) briefly to some of your points: We will reference a logic textbook as suggested; our
41 (overpowered) proof system generalizes standard Hilbert systems and is thus indeed complete. Yes, "satisfies" should
42 have been "consistent with". Yes, the pairwise equivalence should receive more space in a user-facing explanation;
43 we shortened this too much due to page limit. Agreed, we will avoid using " \forall " in the metalanguage. The standard
44 homogeneity axiom is equivalent to the conjunction of MULT and SIMP; we introduced new names since we handle
45 the two parts separately. Thanks for pointing out the resemblance to structures from compilation complexity; we will
46 reference and think more deeply about this. We agree that the tops-only axiom for plurality is unappealing. There
47 are some nicer characterizations of plurality using independence of Pareto-dominated alternatives, but extending our
48 lower-bound technique to the variable-agenda setting will require more work. Your proposed axiomatization of AV is
49 great, and we think it fits into our framework! In particular, along the lines you sketch, it should be possible to derive
50 instances of the cancellation axiom from BASIC-TIE and COMBINE.