

# A theory of measuring, electing, and ranking

Michel Balinski<sup>a</sup> and Rida Laraki

École Polytechnique and Centre National de la Recherche Scientifique, 1 rue Descartes, 75005 Paris, France

Communicated by Ralph E. Gomory, Alfred P. Sloan Foundation, New York, NY, March 26, 2007 (received for review October 27, 2006)

**The impossibility theorems that abound in the theory of social choice show that there can be no satisfactory method for electing and ranking in the context of the traditional, 700-year-old model. A more realistic model, whose antecedents may be traced to Laplace and Galton, leads to a new theory that avoids all impossibilities with a simple and eminently practical method, “the majority judgement.” It has already been tested.**

The theory of social choice concerns methods for amalgamating the appreciations or evaluations of many individuals into one collective appreciation or evaluation. It has two principal applications. (i) Voting: electors in a democracy choose one among several candidates, or committee members decide on one among several courses of action. (ii) Jury decisions: judges evaluate competitors (e.g., figure skaters, gymnasts, pianists, wines, etc.) and rank them or classify them by level of excellence.<sup>b</sup>

The fundamental problem is to find a social decision function (SDF) whose inputs are messages of judges or voters and whose outputs are the jury or electoral decisions, usually rank-orderings of competitors and winners. Much of the theory of social choice has blurred the distinction between a judge’s complex aims, ends, purposes and wishes, in short, his or her preferences or utilities, and the messages he or she is allowed to send.<sup>c</sup>

In the traditional model, consecrated by some seven centuries of use,<sup>d</sup> each individual judge’s or voter’s rank ordering of the competitors is at once his or her message and his or her preferences. Does it mean that the judge prefers this rank ordering above all others; or, that the judge wishes the first competitor on his list to be the winner, the second to be the winner if the first is not, the third to be the winner if the first two are not, and so on down the list; or, is the rank ordering required and chosen strategically by the judge given his or her “true” rank ordering.

In the real world, a judge’s message is simply a message, nothing more. It depends on the judge’s preferences, but it is not and cannot be his or her preferences. In the real world, a judge’s or a voter’s preferences or utilities depends on a host of factors that include the decision (or output), the messages of the other judges (a judge or voter may wish to differ from the others, or on the contrary resemble the others), the social decision function that is used (a judge may prefer a decision given by “democratic” function to one rendered by an “oligarchique” function, or the contrary), and the message he or she thinks is the right one (a judge may prefer honest behavior, or not).

Kenneth Arrow (5), in the first deep theoretical analysis of the theory of social choice, uses the traditional model: each judge’s input message is a rank ordering, routinely interpreted to be a complete expression of his “preferences” (strategic considerations are absent); the output is a rank ordering and a winner (the first-ranked competitor of the order). His celebrated “impossibility” theorem shows that there exists no social welfare function (SWF) satisfying three reasonable properties for obtaining a decision given any inputs (unless there are only two competitors). Amartya Sen (6) models each judge’s inputs as a numerical “utility” over the competitors, i.e., the judge assigns a real number to every competitor; the output is a rank ordering whose utility to a judge is not specified. The model has theoretical interest but no practical significance because a voter’s individual utility is a much more complex concept. In any case, Arrow’s theorem emerges again unless the utilities are assumed to be comparable (that bugbear of

economists!). The model used to prove the well known Gibbard–Satterthwaite (7, 8) impossibility theorem assumes the output is a winner (indeed, how could preferences be modelled if the output were a rank-ordering?); each judge has “true” preferences expressed as a rank ordering; but a judge’s input is a strategically chosen rank ordering. The theorem states that there exists no social choice function that makes it a dominant strategy for every judge to report his true preferences.

Refined, extended, and reformulated in many variants, the traditional approach has continued to produce a host of related impossibility theorems. We add to this list a negative theorem of a new kind: a fundamental incompatibility between winners and rank orderings as outputs of the traditional model. It devolves from a simple observation: if the output is to be a rank ordering and inputs are interpreted as preferences, should not an individual’s input message be his preferences over rank orderings rather than a single rank ordering?

Given all of these negative results, it is not surprising that the debate over what method of voting should be used in practice goes on unabated. By and large, it may be said to pit the supporters of Lull (alias Condorcet) against those of Cusanus (alias Borda), though some argue for a new method, “approval voting” (9), and diverse hybrids are regularly proposed.

We contend that (i) Arrow’s and all the other impossibility and incompatibility results show that the fundamental problem has no acceptable solution in the context of the traditional model. (ii) The traditional approach does not adequately model the messages or the purposes of the judges and voters. (iii) A new model is necessary.

Practice, curiously enough, suggests a different formulation of the inputs. Olympic competitions in figure skating and gymnastics, wine competitions, competitions among pianists, flautists, or orchestras, etc., all use measures or grades. As Lord Kelvin proclaimed, “If you cannot measure, your knowledge is meager and unsatisfactory.” Indeed, Arrow (5) himself states “there are essentially two methods by which social choices can be made, voting, . . . and the market mechanism”; the second uses a measure: price expressed in terms of money.

A measure or grade is a message that has strictly nothing to do with a utility. A judge may dislike a wine and yet give it a high grade because of its merits; he or she may also like a wine and yet, with great satisfaction, give it a low grade because of its demerits. A measure provides a common language, be it numerical, ordinal or

Author contributions: M.B. and R.L. designed research, performed research, and wrote the paper.

The authors declare no conflict of interest.

Abbreviations: SWF, social welfare function; SGF, social grading function.

<sup>a</sup>To whom correspondence should be addressed. E-mail: michel.balinski@shs.polytechnique.fr.

<sup>b</sup>Our thesis is capsuled in this article. A complete account, including proofs, many other results, and references, is given in a forthcoming book: *One-Value, One-Vote: Measuring, Electing and Ranking*.

<sup>c</sup>The word “preferences” misleads: voters do not merely express what they prefer, they may well express what they believe is right (1); a judge in a court of justice is supposed to evaluate conformity with the law, not merely express his preferences. In fact, the real, deep preferences of a judge or voter is a complicated function that depends on the SDF itself.

<sup>d</sup>Ramon Lull proposed a refinement of Condorcet’s method in 1299: it is known today as Copeland’s method. Nicolaus Cusanus put forth what is today known as Borda’s method in 1433 (2–4).

© 2007 by The National Academy of Sciences of the USA

verbal, to grade and classify. In this respect, Arrow's theorem means that, without a common language, there can be no consistent collective decision.

When the messages are grades expressed in a common language, then one method of classifying competitors, candidates, or alternatives, the majority-grade, and one method for ranking them, the majority-ranking, emerge as the only ones that satisfy each of various desirable properties. They are compatible. Moreover, they best resist strategic manipulations of judges and voters under varying assumptions concerning the judges' and voters' preferences or utilities.

### The Traditional Model

There are a finite set of *competitors* (alternatives, candidates, performances, or competing goods)  $C = \{A, \dots, I, \dots, Z\}$  and a finite set of  $n$  *judges* (or voters)  $J = \{1, \dots, j, \dots, n\}$ . Each judge's input *message* is a rank ordering of the competitors. Together, all of the input messages constitute a *preference profile* (in keeping with traditional terminology we use the word "preference(s)" as a synonym for rank orderings in this section). A SWF renders an output—a rank ordering—for any inputs or preference profile.

Take  $A > B$  to mean that a judge ranks  $A$  ahead of  $B$ , and  $A >_S B$  to mean that the SWF (or "Society") ranks  $A$  ahead of  $B$ ; in examples an integer followed by a rank ordering is the number of judges sending that input message.

Condorcet is the first to have realized the essential difficulty of the problem. Consider one of his examples with the 60-judge preference profile:

$$23: A > B > C \quad 2: B > A > C \quad 17: B > C > A \\ 10: C > A > B \quad 8: C > B > A.$$

If majority rule decides the order between each pair of competitors separately, the result is  $A >_S B >_S C >_S A$ . This is the Condorcet paradox: no competitor is favored to all others.

Arrow showed that this is an inescapable conundrum. He imposed three conditions that any SWF should satisfy. (i) *Unanimity*: If every judge prefers competitor  $A$  to competitor  $B$ , i.e.,  $A > B$ , then the SWF ranks  $A$  ahead of  $B$ , i.e.,  $A >_S B$ . (ii) *Non-dictatorship*: The input of no one judge can determine the output of the SWF whatever the inputs of all the other judges. (iii) *Independence of irrelevant alternatives* (IIA): whether the SWF yields  $A >_S B$  or the contrary  $A <_S B$  depends only on the judges' preferences between  $A$  and  $B$ . His theorem shows that, when there are at least three competitors, there is no SWF that satisfies the three conditions for all possible inputs of the judges.

Nevertheless, people vote and judges rank, but how? A judge accords  $k$  *Borda-points* to a competitor if  $k$  opponents are ranked below him.<sup>c</sup> A competitor's *Borda-score* is the sum of his Borda points over all judges; equivalently, it is the sum of the votes he receives in all pair by pair votes (11). The *Borda-ranking* ranks the competitors by their Borda-scores, from highest to lowest; the highest designates the *Borda-winner*. The Borda-ranking for Condorcet's 60-judge example is (each competitor's Borda-score is in parentheses):  $B(69) >_S A(58) >_S C(53)$ .

Condorcet attacked Borda's method. His argument was that, when there exists a *Condorcet-winner*, a competitor who has a majority against every other competitor, then he must be the winner, a property that Borda's method violates, as the following 81-judge example of Condorcet shows:

$$30: A > B > C \quad 1: A > C > B \quad 29: B > A > C \\ 10: B > C > A \quad 10: C > A > B \quad 1: C > B > A.$$

<sup>c</sup>Laplace (10) justified the Borda-points by imagining that each judge wishes to assign a positive real score in some interval  $[0, R]$  to each competitor but is asked instead to rank them. Laplace computed the average of the lowest points, of the next to lowest, on up to the highest, and found them to be proportional to the Borda points.

Here the Borda-ranking is  $B(109) >_S A(101) >_S C(33)$ , yet  $A$  is the Condorcet-winner.

Suppose a profile is split into two parts and a method is applied to each to obtain solutions  $S_1$  and  $S_2$  that have an element in common. The method is *join-consistent* if it selects a solution from  $S_1 \cap S_2$  for the entire profile. Young (12) introduced this idea to characterize Borda's method (and positional methods). Saari (13) reinvented it and applied it to Condorcet's example. Let the first part of the profile be

$$20: A > B > C \quad 28: B > A > C,$$

and the second part be

$$10: A > B > C \quad 10: B > C > A \quad 10: C > A > B. \\ 1: A > C > B \quad 1: C > B > A \quad 1: B > A > C.$$

Each line of the second part is a *Condorcet-component*, a perfect symmetry among the candidates.  $B$  is the clear winner of the first; symmetry shows that all the candidates are tied in the second part; join-consistency implies  $B$  is the winner for the entire profile. This result indicates that the Condorcet-winner, when he exists, is certainly not the candidate who should win in every case! Borda's method (and any "positional" method) avoids this difficulty because the total number of points awarded to every competitor in a Condorcet-component is the same.

Saari (14) then asserted that "all election difficulties" come from Condorcet-components and, to a lesser extent, from more intricate symmetries in the preference profiles; and concluded, "the [Borda-count] applied to all  $n$ -candidates is the unique ranking which avoids all of the indicated problems."<sup>f</sup> However, is Borda's method good for ranking, or for designating winners, or both?

Condorcet (15) proposed a method explicitly for ranking (as was recognized by Young, ref. 16). A voter contributes  $k$  Condorcet-points to an arbitrary rank ordering  $A >_S B >_S C >_S \dots >_S Z$  if his input agrees in  $k$  pair-by-pair comparisons. The *Condorcet-count* of the rank ordering is the sum of the *Condorcet-points* over all voters. The *Condorcet-ranking* is the ranking that maximizes the Condorcet-count. It ranks the Condorcet winner first, and the *Condorcet-loser*, a competitor who loses against every other competitor, last, when either exist. Is it good for ranking, or for designating winners and losers, or both?

The two outputs, a rank-ordering and a winner, have usually been treated as two sides of one coin: given a rank-ordering, the winner is the first-placed competitor; given a mechanism for determining a winner, place him first then apply the mechanism to the remaining competitors and place that winner second, and continue. These are questionable practices.

To see why, consider the preference profile

$$333: A > B > C \quad 333: B > C > A \quad 333: C > A > B.$$

The 999 judges constitute a Condorcet-component, and so cancel each other out. Borda and Condorcet agree on the winners:  $A, B$ , and  $C$  are tied. Condorcet (reasonably) says the three stated rank-orderings are tied for first; Borda (ridiculously), says that all six possible rank orderings are tied for first.

Now, consider the same situation with one additional judge  $A > C > B$ . Borda (reasonably) declares  $A$  the winner and  $B$  the loser, but (ridiculously) the ranking  $A >_S C >_S B$ , because only one judge agrees with it, 666 partially agree, and 333 totally disagree. Condorcet (reasonably) declares the rankings  $A >_S B >_S C$  and  $C >_S A >_S B$  are tied: 333 agree and 667 partially agree. However,  $A$  and  $C$  should certainly not be tied as winners.

<sup>f</sup>Saari (13) proposes "Instant-Borda-Runoff" to counter manipulation: namely, obtain the Borda ranking, drop the bottom candidate, and repeat until one candidate remains. This method always elects the Condorcet winner when he exists.

Borda's method is appropriate for designating winners and losers, Condorcet's method is appropriate for designating rank orderings, a fact already appreciated by Young (17). In fact, the situation is much worse: the two outputs cannot be reconciled.

A SWF is *winner-loser-unanimous* if, whenever all voters rank a candidate first (respectively, last), he is the winner (the loser). It is *choice-compatible* if whenever all voters rank a candidate first (respectively, last) and a Condorcet-component is added to the profile, that candidate is the winner (the loser). It is *rank-compatible* if, whenever a winner is removed from the set of candidates, the new ranking on the remaining candidates agrees with original ranking. Borda's method is choice-compatible but not rank-compatible; Condorcet's is rank-compatible but not choice-compatible.

**Theorem 1 (Incompatibility).** *There exists no winner-loser-unanimous, choice- and rank-compatible method.*

**Grading: The Basic Model**

A thorough investigation of practice shows that scores, measures, or grades have been invented to classify and to rank in an incredibly wide variety of circumstances. Practical people needing practical solutions have increasingly devised mechanisms to transform judges' grades (instead of rank orderings) into a jury's grades to determine final rank orderings. A set of grades (e.g., numbers from 0 to 20, to 25, or to 100; medal nominations from none to bronze, silver or gold; letters from *F* to *A*; or words or phrases from bad to excellent) becomes, in effect, a common language used to assess performances, just as grades determine the standing of students in schools and universities.

Formally, a *common language*  $\Lambda$  is a set of grades  $\alpha, \beta$ , etc., that are strictly ordered. It may be finite or an interval of the real numbers.  $\alpha \geq \beta$  means that either  $\alpha$  is a higher grade than  $\beta$ , in symbols,  $\alpha > \beta$ , or  $\alpha = \beta$ .

There are a finite set of  $m$  competitors  $C = \{A, \dots, I, \dots, Z\}$  and a finite set of  $n$  judges (or voters)  $J = \{1, \dots, j, \dots, n\}$ . A problem is completely specified by an *input* or *profile*  $\Phi = \Phi(C, J)$ : an  $m$  by  $n$  matrix of the grades  $\Phi(I, j) \in \Lambda$  assigned by each of the judges  $j \in J$  to each of the competitors  $I \in C$ .

A *method of grading* is a function  $F$  that assigns to any input or profile  $\Phi$  one *output* or *final grade* in the same language for every competitor  $F: \Lambda^{m \times n} \rightarrow \Lambda^m$ . Designed to assign grades, it must satisfy certain basic properties.

**Axiom 1.**  *$F$  is neutral:  $F(\rho\Phi) = \rho F(\Phi)$ , for any permutation  $\rho$  of the competitors (or rows).*

**Axiom 2.**  *$F$  is anonymous:  $F(\Phi\tau) = F(\Phi)$ , for any permutation  $\tau$  of the judges (or columns).*

**Axiom 3.**  *$F$  is unanimous: If a competitor is given an identical grade  $\alpha$  by every judge, then  $F$  assigns him the grade  $\alpha$ .*

**Axiom 4.**  *$F$  is monotonic: If  $\Phi = \Phi'$  except that one or more judges give higher grades to competitor  $I$  in  $\Phi$  than in  $\Phi'$ , then  $F(\Phi)(I) \geq F(\Phi')(I)$ . Moreover, if all the judges give higher grades to competitor  $I$  in  $\Phi$  than in  $\Phi'$ , then  $F(\Phi)(I) > F(\Phi')(I)$ .*

**Axiom 5.**  *$F$  is independent of irrelevant alternatives (IIA): if the grades assigned by the judges to a competitor  $I \in C$  in two profiles  $\Phi$  and  $\Phi'$  are the same, then  $F(\Phi)(I) = F(\Phi')(I)$ .*

A function  $f: \Lambda^n \rightarrow \Lambda$  that transforms a set of judge's grades into a single grade will be called an aggregation function if it satisfies the following three properties:

- *anonymity:*  $f(\dots, \alpha, \dots, \beta, \dots) = f(\dots, \beta, \dots, \alpha, \dots)$ ;
- *unanimity:*  $f(\alpha, \alpha, \dots, \alpha) = \alpha$ ; and
- *monotonicity:*

$$\alpha_j \leq \beta_j \Rightarrow f(\alpha_1, \dots, \alpha_j, \dots, \alpha_n) \leq f(\alpha_1, \dots, \beta_j, \dots, \alpha_n)$$

and

$$\alpha_1 < \beta_1, \dots, \alpha_n < \beta_n \Rightarrow f(\alpha_1, \dots, \alpha_n) < f(\beta_1, \dots, \beta_n).$$

**Theorem 2 (Possibility).** *A method of grading  $F$  satisfies the five axioms if and only if  $F(\Phi)(I) = f(\Phi(I))$  for every  $I \in C$ , for some one aggregation function  $f$ .*

The average or mean value function is the universally used aggregation function in practice, though sometimes highest and lowest grades are dropped. This means that the output language is almost always richer than the language of the input grades (inputs are usually restricted to discrete levels).

In conformity with most practical applications, the common language is parametrized as a subset of real numbers and whatever aggregation is used, small changes in the parametrization or the input grades should imply small changes in the output or the final grades. Hence, even if the initial language is finite, all possible parametrizations must be considered. It is thus natural to take the common language to be  $[0, R]$  for some positive real  $R$  (as did Laplace), and impose:

**Axiom 6.**  *$F$  (and its aggregation function  $f$ ) is continuous.*

A *social grading function* (SGF)  $F$  is a method of grading that satisfies the six axioms of the basic model.

Thus,  $F$  defines, and is defined by, a unique continuous aggregation function  $f$ . In the sequel, because a SGF and its aggregation function go hand in hand, properties are defined in terms of aggregation functions, theorems stated in terms of SGFs. Also,  $r = (r_1, \dots, r_n)$  represents a competitor's grades, superscripts designate competitors.

Enriching a language by embedding it into a real interval opens the door to many more methods of grading, but it will turn out that the aggregation functions that emerge as those that must be used are directly applicable in the seemingly more restrictive finite languages as well.

**Order Functions**

A judge of the jury knows the SGF (equivalently, the aggregation function  $f$ ) that determines the final grades: what strategies will he use in the "game" of assigning his grades? A judge undoubtedly wishes to give the grade he believes is the "right one;" he may, however, assign it so that the final grade is as close as possible to the "right one;" or, he may try to manipulate the outcome for extraneous reasons (as did a judge of the pairs figure skating in the 2002 Olympic games). This is why, in practice, highest (one or two) and lowest (one or two) grades are often eliminated.

The "utility" of a judge  $j$  is a complicated function  $u_j(\mathbf{r}^*, \mathbf{r}, f, \Lambda)$ , where  $\mathbf{r}^* = (r_1^*, \dots, r_n^*)$  are the grades the judges believe are the right ones and  $\mathbf{r} = (r_1, \dots, r_n)$ , the grades they give. The utility of judge  $j$  might include a term  $-|r_j^* - r_j|$  if he wished to grade honestly; it might contain a term  $-\sum_{i \neq j} |r_i^* - r_i|$  if he wished that the other judges graded honestly; it might include a term  $-|\Lambda - \Lambda_j^*|$  if he wished a language  $\Lambda_j^*$  were used; and it is often assumed to be "single-peaked,"  $u_j(\mathbf{r}^*, \mathbf{r}, f, \Lambda) = -|r_j^* - f(r_1, \dots, r_n)|$ . In fact, judges' utilities, judges' beliefs, their beliefs about the others' beliefs, etc., are all completely unknown and change from one competition to another. The methods we develop depend only on what in practice can be known, as does Vickrey's second price auction mechanism. So (unlike "mechanism design") judges' utilities are never explicitly assumed. The methods that are singled out are nevertheless "strategy-proof" for large classes of reasonable "utilities;" when they are not, they best combat manipulability.

Suppose that  $r$  is a competitor's final grade. An aggregation function is *strategy-proof-in-grading* if, when a judge's input grade is  $r^+ > r$ , any change in his input can only lead to a lower grade; and



if, when a judge's input grade is  $r^- < r$ , any change in his input, can only lead to a higher grade.

Strategy-proof-in-grading implies that it is a *dominant strategy* for a judge to honestly assign the grade he believes is the correct one, whenever the more a final grade deviates from the correct one the less he likes it ("single-peaked preference," a reasonable assumption for most judges who grade). There is a class of SGFs that is easily seen to be strategy-proof: the order functions.

The  $k$ th highest grade is called the  *$k$ th-order function*  $f^k$ .

**Theorem 3.** *The unique strategy-proof-in-grading SGFs are the order functions.*

They are "group strategy-proof-in-grading" as well. [Moulin (18) proves a related technical result but in an entirely different context.] The result holds without the continuity assumption and also when the language is finite.

How can the effects of strategic manipulation be countered when judges' appreciations or utilities are more complex? To manipulate, a judge must be able to raise or to lower the final grade by raising or lowering the grade he assigns. In some situations, a judge can only change the final grade by increasing his grade, in others only by decreasing his grade. A judge who can both lower and raise the final grade has greater opportunity to manipulate.

**Theorem 4.** *There exists no SGF that, for every profile of grades, prevents every judge from both increasing and decreasing the final grade. The unique SGFs for which at most one judge may both increase and decrease a final grade are the order functions.*

Given an aggregation function  $f$  and input grades  $\mathbf{r} = (r_1, \dots, r_n)$ , let  $\mu^-(f, \mathbf{r})$  be number of judges who can decrease the final grade,  $\mu^+(f, \mathbf{r})$  be the number of judges who can increase the final grade, and  $\mu(f, \mathbf{r}) = \mu^-(f, \mathbf{r}) + \mu^+(f, \mathbf{r})$ . Define the *manipulability* of  $f$ , to be

$$\mu(f) = \max_{\mathbf{r}=(r_1, \dots, r_n)} \mu(f, \mathbf{r}).$$

At worst, a judge can both increase and decrease the final grade, so  $\mu(f) \leq 2n$ . In particular, when  $f$  is taken to be the arithmetic mean of the grades (as does Borda's method) the manipulability is maximized,  $\mu(f) = 2n$ . On the other hand, when  $f$  is the  $k$ th-order function,  $\mu(f) = n + 1$ .

**Theorem 5.** *The unique SGFs that minimize manipulability are the order functions.*

Suppose that, after the members of a jury have assigned their grades, some judge wishes to revise his grade by assigning a grade closer to the final grade of the jury: more consensus for that final grade should confirm it. An aggregation function  $f$  is *reinforcing* when  $f(r_1, \dots, r_k, \dots, r_n) = r$  and  $r_k > \hat{r}_k \geq r$  or  $r \geq \hat{r}_k > r_k$  implies  $f(r_1, \dots, \hat{r}_k, \dots, r_n) = r$ .

**Theorem 6.** *The unique reinforcing SGFs are the order functions.*

If every judge assigns a grade in a subset of the grades, then the final grade should belong to that subset. This may be seen as restricting outputs to the language of inputs, or generalizing unanimity. An aggregation function  $f$  *conforms* with the assigned grades if  $\{r_1, \dots, r_n\} \subset S$  implies  $f(r_1, \dots, r_n) \in S$ .

**Theorem 7.** *The unique SGFs that conform with the assigned grades are the order functions.*

The particular language used in grading should make no difference in the ultimate outcomes. An aggregation function should give equivalent grades when one language is faithfully translated into another. This is the "meaningfulness" problem of measurement theory (19) in the context of a jury decision. An aggregation function  $f$  is *language-consistent* if  $f(\phi(r_1), \dots, \phi(r_n)) = \phi(f(r_1, \dots, r_n))$  for all increasing, continuous functions  $\phi : [0, R] \rightarrow [R, \bar{R}]$ ,  $\phi(0) = \underline{R}$ ,  $\phi(R) = \bar{R}$ .

$\dots, r_n))$  for all increasing, continuous functions  $\phi : [0, R] \rightarrow [R, \bar{R}]$ ,  $\phi(0) = \underline{R}$ ,  $\phi(R) = \bar{R}$ .

**Theorem 8.** *The unique language-consistent SGFs are the order functions.*

This result depends crucially on the judges using a common language. When there is no common language, a judge's only meaningful input is the order of his grades. An aggregation function  $f$  is *preference-consistent* if  $f(r_1, \dots, r_n) \geq f(s_1, \dots, s_n)$  implies  $f(\phi_1(r_1), \dots, \phi_n(r_n)) \geq f(\phi_1(s_1), \dots, \phi_n(s_n))$ , for all increasing, continuous functions  $\phi_j : [0, R] \rightarrow [R, \bar{R}]$ ,  $\phi_j(0) = \underline{R}$ ,  $\phi_j(R) = \bar{R}$ .

**Theorem 9. (Arrow's Impossibility).** *There exists no preference-consistent SGF.*

This theorem shows that, to arrive at meaningful final grades, it is essential for judges to share a common language: otherwise, the road is barred by Arrow's fundamental result. However, that only stands to reason: imagine the leaders of the world's powers negotiating an agreement with no common language (and no translators)!<sup>9</sup>

### The Majority-Grade

The evidence supports the use of order functions when juries grade. There are many such functions. Different arguments single out one. Sir Francis Galton (23) had the key idea just one century ago, namely: "[The] middlemost estimate, the number of votes that it is too high being exactly balanced by the number of votes that it is too low. Every other estimate is condemned by a majority of voters as being either too high or too low . . . The number of voters may be odd or even. If odd, there is one middlemost value . . . If the number of voters be even, there are two middlemost values, the mean of which must be taken . . ." He erred in the even case.

A *middlemost* aggregation function  $f$ , for  $r_1 \geq \dots \geq r_n$ , is

$$f(r_1, \dots, r_n) = r_{(n+1)/2} \text{ when } n \text{ is odd, and} \\ r_{n/2} \geq f(r_1, \dots, r_n) \geq r_{(n+2)/2} \text{ when } n \text{ is even.}$$

When  $n$  is odd, it is the order function  $f^{(n+1)/2}$ . When  $n$  is even, there are infinitely many; in particular,  $f^{n/2}$  and  $f^{(n+2)/2}$  are the upper-middlemost and lower-middlemost order functions. The middlemost interval is the  $r_{(n+1)/2}$  when  $n$  is odd, and  $[r_{(n+2)/2}, r_{n/2}]$  when  $n$  is even.

Whatever the parity of  $n$ , every grade other than a grade in the middlemost interval is condemned by an absolute majority of the judges as being either too high or too low.

**Theorem 10.** *The unique aggregation functions that assign a final grade of  $r$  when a majority of judges assign  $r$  are the middlemost.*

Practical mechanisms of grading often eliminate extremes to counter cheaters, to guard against cranks, and to emphasize the significance of place in order rather than magnitude. A SGF *counters crankiness*<sup>h</sup> if for  $r_1 \geq \dots \geq r_n$ ,  $n \geq 3$ , its aggregation function  $f$  satisfies  $f(r_1, r_2, \dots, r_{n-1}, r_n) = f(r_2, \dots, r_{n-1})$ , where in going from left to right the highest and lowest grades have been dropped (the two  $f$ s are, in fact, different, but expressing the idea in this manner simplifies notation). Iterating,  $f(r_1, r_2, \dots, r_{n-1}, r_n) = f(r_+, r_-)$  where  $[r_-, r_+]$  is the middlemost interval (a point  $r_- = r_+$  when  $n$  is odd).

When a judge dislikes a final grade the further it departs from his ideal grade, it is a dominant strategy for him to assign his ideal

<sup>9</sup>Properties close to language- and preference-consistency are known under different names in the literature on measurement theory; in particular, Theorems 8 and 9 are known in one guise or another (20–21). Work on welfarism (22) initiated by Sen (6) has considered similar invariance properties.

<sup>h</sup>The word honors Galton (23), who wished to avoid giving "power to 'cranks' in proportion to their crankiness."

grade. But judges may have different incentives. A judge may wish to either increase or decrease the final grade. The  $k$ th-order function allows  $n - k + 1$  judges to increase the final grade and  $k$  to decrease it. It is desirable to thwart potential manipulation as much as possible. Letting  $\lambda$  be the probability a judge wishes to increase the grade and  $1 - \lambda$  that he wishes to decrease it, the probability of effective-manipulability of the aggregation function  $f$  is

$$EM(f) = \max_{\mathbf{r}=(r_1, \dots, r_n)} \max_{0 \leq \lambda \leq 1} \frac{\lambda \mu^+(f, \mathbf{r}) + (1 - \lambda) \mu^-(f, \mathbf{r})}{n}$$

**Theorem 11.** *The unique aggregation functions that minimize the probability of effective-manipulability or that counter crankiness are the middlemost that depend only on the middlemost interval.*

Many physical measures have the property that equal intervals have the same significance: they are “interval measures” in the jargon of measurement theory. The grades invented to assign to competing skaters, pianists or politicians could be interval measures, but more likely are not. As a grade approaches “perfection,” each additional point often represents much more than an additional point added to a middling grade; and at the other end of the scale, the same phenomenon exists. It is reasonable to suppose that an interval measure exists in theory. In fact, points are routinely added and averaged, so treated as if they were interval measures. This is why some parametrizations attempt to linearize the language. It suffices to postulate the existence of much less to imply the existence of an interval measure.

For example, suppose there exists a distance function  $d$  that measures the judge’s discontent: when the judge assigns the grade  $r$  and the final grade is  $s$ , his disutility is  $d(r, s) \geq 0$ . Thus,  $d$  satisfies:  $d(r, r) = 0$ ,  $d(r, s) = d(s, r)$ , and  $r < s < t$  implies  $d(r, s) + d(s, t) = d(r, t)$ . The last equation says that the improvement in a competitor’s performance in going from a grade of  $r$  to  $s$  plus the improvement in going from  $s$  to  $t$  equals that of going from  $r$  to  $t$ ; or, that the disutility of a judge who believes the grade should be  $r$  when the final grade is  $t$  equals his disutility when the final grade is  $s$  plus his disutility when he believes it should be  $s$  when the final grade is  $t$ . This accommodates the possibility that, for example, on a scale of 0 to 100,  $d(98, 99) = 5d(75, 76)$ . Several arguments suggest that it is reasonable to assume that all the judges view changes in performances similarly: one judge teaches all others; or, the rules impose by fiat that they do; or, equity among the judges demands that their disutilities must be modelled identically.

A SGF with aggregation function  $f$  maximizes welfare when the final grade  $f(r_1, \dots, r_n) = r$  minimizes the total disutility of all of the judges,  $\Delta(r) = \sum_{j \in \mathcal{J}} d(r, r_j)$ .

**Theorem 12.** *The unique aggregation functions that maximize welfare are the middlemost.*

Thus imposing an equity condition, namely, that judges compare performances with the same measure, together with the assumption that the measure is a distance function, implies that the optimal mechanism must be a majority decision. A distance function is equivalent to the existence of an interval measure (not necessarily compact).<sup>1</sup>

**Characterization:** A SGF rewards consensus when all of  $A$ ’s grades strictly belong to the middlemost interval of  $B$ ’s grades implies that  $A$ ’s final grade is higher than  $B$ ’s final grade.

The majority-grade  $f^{maj}$  is the SGF defined by the order function  $f^{(n+1)/2}$  when  $n$  is odd, and by the lower-middlemost order function  $f^{(n+2)/2}$  when  $n$  is even.

**Theorem 13.** *The unique middlemost aggregation function that rewards consensus is the majority-grade  $f^{maj}$ .*

<sup>1</sup>Defining  $\phi(s) = d(R/2, s)$  if  $s \geq R/2$  and  $-\phi(s) = d(s, R/2)$  if  $s < R/2$  implies  $d(r, s) = |\phi(r) - \phi(s)|$ .

## The Majority-Ranking

A competitor bestowed a higher grade than another is naturally ranked higher in the order of the competitors than the other: grades imply orders. The essential incompatibility between the designation of winners (or losers) and rank orderings inherent to the traditional model (Theorem 1) simply does not arise in the context of grades. On the other hand, although in some applications a complete ordering is not sought, e.g., wine competitions, there are other applications, notably sports and elections, where an ordered list from first to last and a clear winner is absolutely necessary.

When rank orderings are the principal goal instead of grades, the strategic behavior of the judges may change. A SGF is strategy-proof-in-ranking if for any judge  $j$  final grades  $r^A < r^B$  opposed to the judge’s grades  $r_j^A > r_j^B$  implies he can neither decrease  $B$ ’s final grade nor increase  $A$ ’s final grade; and it is partially strategy-proof-in-ranking if any judge in the same situation can decrease  $B$ ’s final grade implies he cannot increase  $A$ ’s and if he can increase  $A$ ’s final grade implies he cannot decrease  $B$ ’s.

**Theorem 14.** *There exists no SGF that is strategy-proof-in-ranking. The unique SGFs that are partially strategy-proof-in-ranking are the order functions.*

When the majority-grades of two competitors  $A$  and  $B$  differ, the one with the higher majority-grade ranks ahead of the other. When the majority-grades of two competitors are equal, no more useful information concerning these two competitors can be drawn from this grade.

The majority-ranking ( $>_{maj}$ ) between two competitors is determined by:

- (1) If  $f^{maj}(A) > f^{maj}(B)$  then  $A >_{maj} B$ .
- (2) If  $f^{maj}(A) = f^{maj}(B)$ , the majority-grade is dropped from the grades of each of the competitors, and the procedure is repeated.

**Theorem 15.** *The majority-ranking always ranks one competitor ahead of another unless the two are assigned an identical set of grades by the judges.*

The first-majority-grade of a competitor is the majority grade of the entire jury; the second-majority-grade is the majority grade of the grades that remain after the first-majority grade has been dropped; the  $i$ th majority grade is the majority grade of the grades that remain after the first  $i - 1$  majority grades have been dropped. A competitor’s majority-value is a vector of  $n$  components that assigns, in order, his first, second, third, . . . ,  $n$ th-majority-grades.

**Theorem 16.**  *$A >_{maj} B$  if and only if  $A$ ’s majority-value is lexicographically higher than  $B$ ’s.*

The majority-value assigns a specific value to each competitor expressed in terms of the common language. It may be transformed into a rational number when the language is finite. For example, if the language has ten grades 0, 3, 5, 6, . . . , 11, 13 (Denmark’s school grades; the absence of certain numbers is an attempt to linearize) and a competitor receives the grades (7, 7, 9, 10, 11), then his majority-value is 9, 07100711. Dividing by 1.01010101 rescales the final values so that the minimum is 0, the maximum 13, and a candidate assigned the same grade  $\alpha$  by all judges has a rescaled majority value of  $\alpha$ .

**Characterization:** Given input grades  $\mathbf{r}^A = (r_1^A, \dots, r_n^A)$ ,  $\mathbf{r}^B = (r_1^B, \dots, r_n^B)$  of two competitors, how should they be ranked? Write  $A >_S B$  to mean  $A$  is ranked ahead of  $B$ , and  $A \geq_S B$  to mean either  $A$  is ahead of  $B$  or they are tied. If  $n \geq 3$ , the highest and the lowest grades are  $\mathbf{r}$ ’s residual grades, its set of center grades is obtained by dropping the residual grades.

A social ranking function (SRF) should satisfy several properties. It should be (i) *monotone*: if  $A \geq_S B$  and one judge raises the grade he gives to  $A$  then  $A >_S B$ . It should be (ii) *decisive for the center grades*: the ranking between  $A$  and  $B$  is the ranking determined by

the center grades unless that ranking is a tie; in that case, the ranking is determined by the residual grades. Finally, it should (iii) reward consensus.

**Theorem 17.** *The majority-ranking is the unique monotone SRF that is decisive for the center and rewards consensus.*<sup>1</sup>

The majority-ranking is IIA in the sense of Arrow: the order between two competitors depends only on their respective grades.

**Remark.** This theory is not “cardinal”: adding grades is meaningless. Nor is it “ordinal”: a voter’s message depends on the particular words of the language that is used, so different common languages with the same number of words may lead to different rankings. And yet, a grade in a common language has an absolute meaning.

## Practice

**Jury Decisions.** The majority grade and ranking are described in the new edition of the French “Bible” on wines (24). They were tested in a wine competition, *Les citadelles du vin*, held June 15–17, 2006, in the Bordeaux region of France. A total of 1,247 different wines competed, 60 judges were organized in juries of five members (sometimes fewer). As usual, judges completed the “Sensorial analysis tasting sheet for wine judging competitions” of the *Organisation Internationale de la Vigne et du Vin* for each wine: 14 attributes of the wine were assigned points (from 0 “bad” to 6 or 8 “excellent”); their sum determined whether the wine was bestowed a gold, silver, bronze, or no medal at all. However, the sum misses the point (25) because it “has difficulty in detecting exceptional wines by overly favoring those that are ‘taste-wise correct.’” Moreover, there is strong evidence showing that judges work “backwards,” they first decide the grade they wish to bestow, then assign points to attributes whose sum yield the grade. The judges preferred to answer “For you, this wine is:” with one of five descriptions that constituted the common language in this experiment: Excellent, Very good, Good, Average, or Mediocre. A preliminary evaluation of the experiment concludes: “The ‘majority-grade’ correctly distinguished . . . the wines, in accordance with the traditional objectives of wine competitions. This system seems better adapted than the [old system] . . . [However], the scale of five levels—97% of the grades were confined to three levels—should be extended” (J. Blouin, personal communication).

**Voting.** The majority ranking’s first-ranked candidate designates the winner of an election. Approval voting (AV) uses a common language of two words, 1 “approve” and 0 “disapprove.” The

majority-ranking with a 0,1-language is the approval voting ranking, so AV is a special case of the majority-ranking. AV has been tested in a variety of settings, notably professional scientific societies. It was also tested (26) in parallel with the first round of the French presidential election of 2002, when 16 candidates presented themselves: voters were clearly happy to be able to express themselves better with AV than casting at most one solitary vote. The arguments for and against AV have been cast in the context of the traditional model and have not addressed the real problem. The only solid results assume “dichotomous preferences” (27, 28) but “like” and “dislike” (or the very different “for” and “against”) is much too limited a language.

Why do electors vote at all, when they hardly expect to determine the outcome (1)? They feel the moral imperative to express themselves: why else do so many cast blank votes? A richer language should encourage greater public participation. Exactly what common language should be used in, say, a presidential election, is not a trivial choice. Perhaps it should be the grading system used in the nation’s educational system: from a low of *F* to a high of *A* in the U.S., from 0 to 20 in France, or 0 to 13 in Denmark. Alternatively, an election of an official might ask each voter: “For you, this candidate is Exceptional, Accomplished, Capable, Average, Limited, or Incompetent to undertake the high responsibilities of [the office].” The method was tested in the 2007 French presidential elections ([www.ceco.polytechnique.fr/jugement-majoritaire.html](http://www.ceco.polytechnique.fr/jugement-majoritaire.html)).

**Common Language.** How to define a common language in general remains an open question, though in many applications (e.g., skating, diving, gymnastics, piano competitions) languages already exist. Different applications naturally call for different common languages. Experimentation will be necessary to define a language; and, as is true of any language, it will alter over time. The *Les citadelles du vin* experiment suggests that judges (and voters?) shun the highest and lowest grades. It may be best to define a language with an even number of words in order to prevent voting in the middle, or not. The nature of the words or numbers used will illicit different voting and judging behavior: the words themselves matter! The environment in which judging and voting take place may also. Just imagine, how would responsible voters behave were they to read Ramon Lull’s (3) solemn proclamation of 1299 before casting their ballots:

. . . [It] is necessary to ascertain that in the election three things should be considered, of which the first is honesty and holiness of life, the second is knowledge and wisdom, and the third is a suitable disposition of the heart. Each person having a vote in the chapter should take an oath by the holy gospels of God to consider these three things and to always elect the person in whom they are best [embodied].

<sup>1</sup>The majority-value with Borda or Condorcet points as inputs provide SWFs for the traditional model that combat strategic manipulation. The first-majority-grade with Borda points was used to rank figure skaters prior to 2004, with *ad hoc* rules to resolve ties.

- Goodin R, Roberts K (1975) *Am Pol Sci Rev* 69:926–928.
- McLean I (1990) *Soc Choice Welfare* 7:99–108.
- Hägele G, Pukelsheim F (2001) *Stud Lulliana* 41:3–38.
- Hägele G, Pukelsheim F (2007) in *The Church, the Councils and Reform: Lessons from the Fifteenth Century*, eds Christianson HG, Izbicki TM, Bellitto CM (Catholic Univ of America Press, Washington, DC), in press.
- Arrow KJ (1951) *Social Choice and Individual Values* (Wiley, New York).
- Sen A (1970) *Collective Choice and Social Welfare* (Holden-Day, San Francisco).
- Gibbard A (1973) *Econometrica* 41:587–601.
- Satterthwaite M (1973) *J Econ Theory* 10:187–217.
- Brams SJ, Fishburn PC (1983) *Approval Voting* (Birkhäuser, Boston).
- de Laplace P-S (1820) in *Œuvres Complètes de Laplace*, tome 7, 3rd Ed, pp v and clii–cliij.
- de Borda J-C (1784) *Histoire de l’Académie Royale des Sciences*, 657–665.
- Young HP (1975) *SIAM J Appl Math* 28:824–838.
- Saari D (2001) *Chaotic Elections! A Mathematician Looks at Voting* (Am Math Soc, Washington, DC), pp 134 and 103.
- Saari D (2000) *Econ Theory* 15:57.
- de Condorcet JAC (1785) *Essai sur l’Application de l’Analyse à la Probabilité des Décisions Rendues à la Pluralité des Voix* (l’Imprimerie Royale, Paris).
- Young HP (1988) *Am Pol Sci Rev* 82:1231–1244.
- Young HP (1986) in *Information Pooling and Group Decision Making*, eds Grofman B, Owen G (JAI, Greenwich, CT), pp 113–122.
- Moulin H (1980) *Public Choice* 35:437–455.
- Krantz DH, Luce RD, Suppes P, Tversky A (1971) *Foundations of Measurement: Additive and Polynomial Representation* (Academic, New York), Vol 1.
- Orlov A (1981) *Math Notes* 30:774–778.
- Kim S-R (1990) *Math Soc Sci* 20:19–36.
- Bossert W, Weymark JA (2004) in *Handbook of Utility Theory: Extensions*, eds Barberà S, Hammond PJ, Seidl C (Kluwer Academic, Boston), Vol 2, pp 1099–1177.
- Galton F (1907) *Nature* 75:414.
- Peynaud E, Blouin J (2006) *Le goût du vin* (Dunod, Paris), pp 104–107.
- Peynaud E, Blouin J (1999) *Découvrir le goût du vin*, (Dunod, Paris), p 109.
- Balinski M, Laraki R, Laslier J-F, van der Straeten K (2002) *Pour la Science*, 13.
- Barberà S, Sonnenschein H, Zhou L (1991) *Econometrica* 59:595–609.
- Bogomolnaia A, Moulin H, Strong R (2005) *J Econ Theory* 122:165–184.